

# MM-LLMs: Recent Advances in MultiModal Large Language Models

Anonymous ACL submission

## Abstract

In the past year, MultiModal Large Language Models (MM-LLMs) have undergone substantial advancements, augmenting off-the-shelf LLMs to support MM inputs or outputs via cost-effective training strategies. The resulting models not only preserve the inherent reasoning and decision-making capabilities of LLMs but also empower a diverse range of MM tasks. In this paper, we provide a comprehensive survey aimed at facilitating further research of MM-LLMs. Initially, we outline general design formulations for model architecture and training pipeline. Subsequently, we introduce a taxonomy encompassing 122 MM-LLMs, each characterized by its specific formulations. Furthermore, we review the performance of selected MM-LLMs on mainstream benchmarks and summarize key training recipes to enhance the potency of MM-LLMs. Finally, we explore promising directions for MM-LLMs while concurrently maintaining a real-time tracking website<sup>1</sup> for the latest developments in the field. We hope that this survey contributes to the ongoing advancement of the MM-LLMs domain.

## 1 Introduction

MultiModal (MM) pre-training research has witnessed significant advancements in recent years, consistently pushing the performance boundaries across a spectrum of downstream tasks (Li et al., 2020; Akbari et al., 2021; Fang et al., 2021; Yan et al., 2021; Li et al., 2021; Radford et al., 2021; Li et al., 2022; Zellers et al., 2022; Zeng et al., 2022b; Yang et al., 2022; Wang et al., 2022a,b). However, as the scale of models and datasets continues to expand, traditional MM models incur substantial computational costs, particularly when trained from scratch. Recognizing that MM research operates at the intersection of various modalities, a logical approach is to capitalize on readily available pre-trained unimodal foundation models, with

<sup>1</sup><https://mm-llms.github.io>

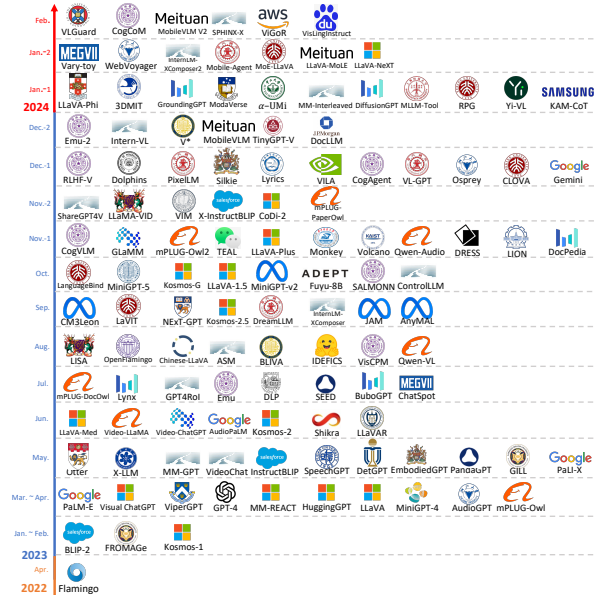


Figure 1: The timeline of MM-LLMs.

a special emphasis on powerful Large Language Models (LLMs) (OpenAI, 2022). This strategy aims to mitigate computational expenses and enhance the efficacy of MM pre-training, leading to the emergence of a novel field: MM-LLMs.

MM-LLMs harness LLMs as the cognitive powerhouse to empower various MM tasks. LLMs contribute desirable properties like robust language generation, zero-shot transfer capabilities, and In-Context Learning (ICL). Concurrently, foundation models in other modalities provide high-quality representations. Considering foundation models from different modalities are individually pre-trained, the core challenge facing MM-LLMs is how to effectively connect LLMs with models in other modalities to enable collaborative inference. The predominant focus within this field has been on refining alignment between modalities and aligning with human intent via a MM Pre-Training (PT) + MM Instruction-Tuning (IT) pipeline.

With the debut of GPT-4(Vision) (OpenAI, 2023)

and Gemini (Team et al., 2023), showcasing impressive MM understanding and generation capabilities, a research fervor on MM-LLMs has been sparked. Initial research primarily focuses on MM content comprehension and text generation, encompassing tasks such as image-text understanding, exemplified by projects like BLIP-2 (Li et al., 2023e), LLaVA (Liu et al., 2023e), MiniGPT-4 (Zhu et al., 2023a), and OpenFlamingo (Awadalla et al., 2023); video-text understanding, as demonstrated by initiatives such as VideoChat (Li et al., 2023f), Video-ChatGPT (Maaz et al., 2023), and LLaMA-VID (Li et al., 2023j); and audio-text understanding, as seen in projects like Qwen-Audio (Chu et al., 2023b). Later, the capabilities of MM-LLMs have been expanded to support specific modality outputs. This includes tasks with image-text output, such as GILL (Koh et al., 2023a), Kosmos-2 (Peng et al., 2023), Emu (Sun et al., 2024), and MiniGPT-5 (Zheng et al., 2023b); as well as speech/audio-text output, exemplified by projects like SpeechGPT (Zhang et al., 2023a) and AudioPaLM (Rubenstein et al., 2023). Recent research endeavors have focused on mimicking human-like any-to-any modality conversion, shedding light on the path to artificial general intelligence. Some efforts aim to amalgamate LLMs with external tools to reach an approaching any-to-any MM comprehension and generation, such as Visual-ChatGPT (Wu et al., 2023a), HuggingGPT (Shen et al., 2023), and AudioGPT (Huang et al., 2023b). Conversely, to mitigate propagated errors in the cascade system, initiatives like NEX-T-GPT (Wu et al., 2023d), CoDi-2 (Tang et al., 2023c), and ModaVerse (Wang et al., 2024c) have developed end-to-end MM-LLMs of arbitrary modalities. The timeline of MM-LLMs is depicted in Figure 1.

In this paper, we present a comprehensive survey aimed at facilitating further research of MM-LLMs. To provide readers with a holistic understanding of MM-LLMs, we initially delineate general design formulations from model architecture (Section 2) and training pipeline (Section 3). We break down the general model architecture into five components: Modality Encoder (Section 2.1), Input Projector (Section 2.2), LLM Backbone (Section 2.3), Output Projector (Section 2.4), and Modality Generator (Section 2.5). The training pipeline elucidates how to enhance a pre-trained text-only LLM to support MM input or output, primarily consisting of two stages: MM PT (Section 3.1) and MM IT (Section 3.2). In that section, we

also provide a summary of mainstream datasets for MM PT and MM IT. Next, we establish a taxonomy encompassing 122 State-of-the-Art (SOTA) MM-LLMs, each characterized by specific formulations, and summarize their development trends in Section 4. In Section 5, we comprehensively review the performance of major MM-LLMs on mainstream benchmarks and distill key training recipes to enhance the efficacy of MM-LLMs. In Section 6, we offer promising directions for MM-LLMs research. Moreover, we have established a website (<https://mm-llms.github.io>) to track the latest progress of MM-LLMs and facilitate crowdsourcing updates. Finally, we summarize the entire paper in Section 7 and discuss related surveys on MM-LLMs in Appendix A. We aspire for our survey to aid researchers in gaining a deeper understanding of this field and to inspire the design of more effective MM-LLMs.

## 2 Model Architecture

In this section, we provide a detailed overview of the five components comprising the general model architecture, along with the implementation choices for each component, as illustrated in Figure 2. MM-LLMs that emphasize MM understanding only include the first three components. During training, Modality Encoder, LLM Backbone, and Modality Generator are generally maintained in a frozen state. The primary optimization emphasis is on Input and Output Projectors. Given that Projectors are lightweight components, the proportion of trainable parameters in MM-LLMs is notably small compared to the total parameter count (typically around 2%). The overall parameter count is contingent on the scale of the core LLM utilized in the MM-LLMs. As a result, MM-LLMs can be efficiently trained to empower various MM tasks.

### 2.1 Modality Encoder

The Modality Encoder (ME) is tasked with encoding inputs from diverse modalities  $I_X$  to obtain corresponding features  $F_X$ , formulated as follows:

$$F_X = \text{ME}_X(I_X). \quad (1)$$

Various pre-trained encoder options  $\text{ME}_X$  exist for handling different modalities, where  $X$  can be image, video, audio, 3D, etc. Next, we will offer a concise introduction organized by modality.

**Visual Modality** For images, there are various optional encoders: **NFNet-F6** (Brock et al.,

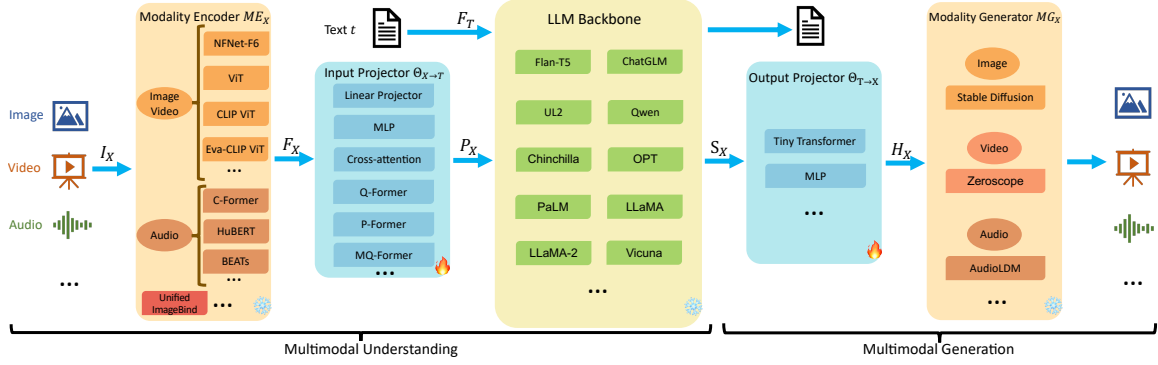


Figure 2: The general model architecture of MM-LLMs and the implementation choices for each component.

2021), **ViT** (Dosovitskiy et al., 2020), **CLIP ViT** (Radford et al., 2021), **Eva-CLIP ViT** (Fang et al., 2023), **BEiT-3** (Wang et al., 2023d), **OpenCLIP** (Cherti et al., 2023), **Grounding-DINO-T** (Zhang et al., 2022b) with Swin-T (Liu et al., 2021b) backbone, **DINOv2** (Oquab et al., 2023), **SAM-HQ** (Kirillov et al., 2023) with MAE (He et al., 2022), **RAM++** (Zhang et al., 2023i) with Swin-B backbone, **InternViT** (Chen et al., 2023j), and **VCoder** (Jain et al., 2023). For videos, they can be uniformly sampled to 5 frames, undergoing the same pre-processing as images.

**Audio Modality** is typically encoded by **C-Former** (Chen et al., 2023b), **HuBERT** (Hsu et al., 2021), **BEATs** (Chen et al., 2023g), **Whisper** (Radford et al., 2023), and **CLAP** (Wu et al., 2023e).

**3D Point Cloud Modality** is typically encoded by **ULIP-2** (Salesforce, 2022) with a PointBERT (Yu et al., 2022) backbone.

Moreover, to handle numerous heterogeneous modal encoders, some MM-LLMs, particularly any-to-any ones, use **ImageBind** (Girdhar et al., 2023), a unified encoder covering six modalities, including image/video, text, audio, heat map, inertial measurement units, and depth. We provide a brief introduction to some mainstream modality encoders in Appendix B.

## 2.2 Input Projector

The Input Projector  $\Theta_{X \rightarrow T}$  is tasked with aligning the encoded features of other modalities  $F_X$  with the text feature space  $T$ . The aligned features as prompts  $P_X$  are then fed into the LLM Backbone alongside the textual features  $F_T$ . Given  $X$ -text dataset  $\{I_X, t\}$ , the goal is to minimize the  $X$ -conditioned text generation loss  $\mathcal{L}_{\text{txt-gen}}$ :

$$\arg \min_{\Theta_{X \rightarrow T}} \mathcal{L}_{\text{txt-gen}}(\text{LLM}(P_X, F_T), t), \quad (2)$$

where  $P_X = \Theta_{X \rightarrow T}(F_X)$ .

The Input Projector can be achieved directly by a **Linear Projector** or Multi-Layer Perceptron (**MLP**), *i.e.*, several linear projectors interleaved with non-linear activation functions. There are also more complex implementations like **Cross-attention**, **Q-Former** (Li et al., 2023e), **P-Former** (Jian et al., 2023), and **MQ-Former** (Lu et al., 2023a). **Cross-attention** (Perceiver Resampler) (Alayrac et al., 2022) uses a set of trainable vectors as queries and the encoded features  $F_X$  as keys to compress the feature sequence to a fixed length. The compressed representation is then fed directly into the LLM or further used for X-Text cross-attention fusion. **Q-Former** extracts relevant features from  $F_X$  with learnable queries, and the selected features are then used as prompts  $P_X$ . Meanwhile, **P-Former** generates "reference prompts", imposing an alignment constraint on the prompts produced by Q-Former. **MQ-Former** conducts a fine-grained alignment of multi-scale visual and textual signals. However, both Q-, P-, MQ-Former require an **additional PT** process for initialization.

## 2.3 LLM Backbone

Taking LLMs (Zhao et al., 2023c; Naveed et al., 2023; Luo et al., 2023) as the core agents, MM-LLMs can inherit some notable properties like zero-shot generalization, few-shot ICL, Chain-of-Thought (CoT), and instruction following. The LLM Backbone processes representations from various modalities, engaging in semantic understanding, reasoning, and decision-making regarding the inputs. It produces (1) direct textual outputs  $t$ , and (2) signal tokens  $S_X$  from other modalities (if any). These signal tokens act as instructions to guide the generator on whether to produce MM contents and,

if affirmative, specifying the content to produce:

$$t, S_X = \text{LLM}(P_X, F_T), \quad (3)$$

where the aligned representations of other modalities  $P_X$  can be considered as soft Prompt-tuning for the LLM. Moreover, some works have introduced Parameter-Efficient Fine-Tuning (PEFT) methods, such as Prefix-tuning (Li and Liang, 2021), LoRA (Hu et al., 2021), and LayerNorm tuning (Zhao et al., 2024). In these cases, the number of additional trainable parameters is exceptionally minimal, even less than 0.1% of the total LLM parameter count. We provide an introduction to mainstream PEFT methods in Appendix C.

The commonly used LLMs in MM-LLMs include **Flan-T5** (Chung et al., 2022), **ChatGLM** (Zeng et al., 2022a), **UL2** (Tay et al., 2022), **Persimmon** (Elsen et al., 2023), **Qwen** (Bai et al., 2023a), **Chinchilla** (Hoffmann et al., 2022), **OPT** (Zhang et al., 2022c), **PaLM** (Chowdhery et al., 2023), **LLaMA** (Touvron et al., 2023a), **LLaMA-2** (Touvron et al., 2023b), and **Vicuna** (Chiang et al., 2023). We provide a brief introduction to some representative LLMs in Appendix D.

## 2.4 Output Projector

The Output Projector  $\Theta_{T \rightarrow X}$  maps the signal token representations  $S_X$  from the LLM Backbone into features  $H_X$  understandable to the following Modality Generator  $\text{MG}_X$ . Given the  $X$ -text dataset  $\{I_X, t\}$ ,  $t$  is first fed into LLM to generate the corresponding  $S_X$ , then mapped into  $H_X$ . To facilitate alignment of the mapped features  $H_X$ , the goal is to minimize the distance between  $H_X$  and the conditional text representations of  $\text{MG}_X$ :

$$\arg \min_{\Theta_{T \rightarrow X}} \mathcal{L}_{\text{mse}}(H_X, \tau_X(t)). \quad (4)$$

The optimization only relies on captioning texts, without utilizing any audio or visual resources  $X$ , where  $H_X = \Theta_{T \rightarrow X}(S_X)$  and  $\tau_X$  is the textual condition encoder in  $\text{MG}_X$ . The Output Projector is implemented by a **Tiny Transformer** with a learnable decoder feature sequence or **MPL**.

## 2.5 Modality Generator

The Modality Generator  $\text{MG}_X$  is tasked with producing outputs in distinct modalities. Commonly, existing works use off-the-shelf Latent Diffusion Models (LDMs) (Song et al., 2021; Bao et al., 2022; Zhao et al., 2022), *i.e.*, **Stable Diffusion** (Rombach

et al., 2022) for image synthesis, **Zeroscope** (Cerspense, 2023) for video synthesis, and **AudioLDM-2** (Liu et al., 2023b,c) for audio synthesis. The features  $H_X$  mapped by the Output Projector serve as conditional inputs in the denoising process to generate MM content. During training, the ground truth content is first transformed into a latent feature  $z_0$  by the pre-trained VAE (Kingma and Welling, 2013). Then, noise  $\epsilon$  is added to  $z_0$  to obtain the noisy latent feature  $z_t$ . A pre-trained Unet (Ronneberger et al., 2015)  $\epsilon_X$  is used to compute the conditional LDM loss  $\mathcal{L}_{X\text{-gen}}$  as follows:

$$\mathcal{L}_{X\text{-gen}} := \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1), t} \|\epsilon - \epsilon_X(z_t, t, H_X)\|_2^2, \quad (5)$$

which optimizes parameters  $\Theta_{X \rightarrow T}$  and  $\Theta_{T \rightarrow X}$  by minimizing  $\mathcal{L}_{X\text{-gen}}$ .

## 3 Training Pipeline

MM-LLMs’ training pipeline can be delineated into two principal stages: MM PT and MM IT.

### 3.1 MM PT

During the PT stage, typically leveraging the X-Text datasets, Input and Output Projectors are trained to achieve alignment among various modalities by optimizing predefined objectives. For MM understanding models, optimization focuses solely on Equation (2), while for MM generation models, optimization involves Equations (2), (4), and (5). In the latter case, Equation (2) also includes the ground-truth signal token sequence.

The X-Text datasets include Image-Text, Video-Text, and Audio-Text, with Image-Text having two types: Image-Text pairs (*e.g.*, **<img1>** <txt1>) and interleaved Image-Text corpus (*e.g.*, <txt1>**<img1>**<txt2><txt3>**<img2>**<txt4>). Details of X-Text datasets are shown in Table 3 of Appendix G.

### 3.2 MM IT

MM IT is a method that entails fine-tuning of pre-trained MM-LLMs using instruction-formatted datasets (Wei et al., 2021). Through this process, MM-LLMs can generalize to unseen tasks by adhering to new instructions, thereby enhancing zero-shot performance. This straightforward yet impactful concept has catalyzed subsequent success in the field of NLP, exemplified by works such as InstructGPT (Ouyang et al., 2022), OPT-IML (Iyer et al., 2022), and InstructBLIP (Dai et al., 2023).

MM IT comprises Supervised Fine-Tuning (SFT) and Reinforcement Learning from Human





Figure 3: Taxonomy for MM-LLMs. I: Image, V: Video, A/S: Audio/Speech, and T: Text. I<sub>B</sub>: Document understanding, I<sub>B</sub>: Output bounding box, I<sub>M</sub>: Output segmentation mask, and I<sub>R</sub>: Output retrieved images.

Feedback (RLHF), aiming to align with human intentions and enhance the interaction capabilities of MM-LLMs. SFT converts part of the PT stage data into an instruction-aware format. Using visual Question-Answer (QA) as an example, various templates may be employed like (1) "**<Image>{Question}**" A short answer to the question is; (2) "**<Image>**" Examine the image and respond to the following question with a brief answer: "**{Question}. Answer:**"; and so on. Next, it fine-tunes pre-trained MM-LLMs using the same optimization objectives. SFT datasets can be structured as either single-turn QA or multi-turn dialogues.

After SFT, RLHF involves further fine-tuning of the model, relying on feedback regarding the MM-LLMs' responses (e.g., Natural Language Feedback (NLF) labeled manually or automatically) (Sun et al., 2023b). This process employs a reinforcement learning algorithm to effectively

integrate the non-differentiable NLF. The model is trained to generate corresponding responses conditioned on the NLF (Chen et al., 2023i; Akyürek et al., 2023). The statistics for SFT and RLHF datasets are presented in Table 4 of Appendix G.

The datasets used by existing MM-LLMs in the MM PT and MM IT stages are diverse, but they are all subsets of the datasets in Tables 3 and 4.

## 4 SOTA MM-LLMs

As shown in Figure 3, we classify the 122 SOTA MM-LLMs from both functional and design perspectives. In the design division, "Tool-using" denotes treating the LLM as black box and providing access to certain MM expert systems to perform specific MM tasks via reasoning, while "End-to-End" signifies that the entire model is trained jointly in an end-to-end manner. Based on the previously defined design formulations, we also

Model	I→O	Modality Encoder	Input Projector	LLM Backbone	Output Projector	Modality Generator	#_PT	#_IT
Flaningo	I+V+T→T	I/V: NFNct-F6	Cross-attention	Chinchilla-1.4B/7B/70B	-	-	-	-
BLIP-2	I+T→T	I: CLIP/Eva-CLIP ViT@224	Q-Former w/ Linear Projector	Flan-T5/OPT	-	-	129M	-
LLaVA	I+T→T	I: CLIP ViT-L/14	Linear Projector	Vicuna-7B/13B	-	-	-	-
MiniGPT-4	I+T→T	I: Eva-CLIP ViT-G/14	Q-Former w/ Linear Projector	Vicuna-13B	-	-	-	-
mPLUG-Owl	I+T→T	I: CLIP ViT-L/14	Cross-attention	LLaMA-7B	-	-	-	-
Otter	I+T→T	I: CLIP ViT-L/14	Cross-attention	LLaMA-7B	-	-	-	-
X-LLM	I+V+A+T→T	I/V: ViT-G; A: C-Former	Q-Former w/ Linear Projector	ChatGLM-6B	-	-	-	-
VideoChat	V+T→T	I: ViT-G	Q-Former w/ Linear Projector	Vicuna	-	-	-	-
InstructBLIP	I+V+T→T	I/V: ViT-G/14@224	Q-Former w/ Linear Projector	Flan-T5/Vicuna	-	-	129M	1.2M
PandaGPT	I+T→T	I: ImageBind	Linear Projector	Vicuna-13B	-	-	-	-
GILL	I+T→T	I: CLIP ViT-L	Linear Projector	OPT-6.7B	Tiny Transformer	I: Stable Diffusion-1.5	-	-
Pali-X	I+T→T	I: ViT	Linear Projector	UL2-32B	-	-	-	-
Video-LLaMA	I+V+A+T→T	I/V: Eva-CLIP ViT-G/14; A: ImageBind	Q-Former w/ Linear Projector	Vicuna/LLaMA	-	-	-	-
Video-ChatGPT	V+T→T	I: CLIP ViT-L/14	Linear Projector	Vicuna-v1.1	-	-	-	-
Shikra	I+T→T+I <sub>h</sub>	I: CLIP ViT-L/14@224	Linear Projector	Vicuna-7B/13B	-	-	600K	5.5M
LLaVAR	I+T→T	I: CLIP ViT-L/14@224 & CLIP ViT-L/14@336	Linear Projector	Vicuna-13B	-	-	-	-
mPLUG-DocOwl	I+T→T	I: CLIP ViT-L/14	Cross-attention	LLaMA-7B	-	-	-	-
Lynx	I+V+T→T	I/V: Eva-CLIP ViT-IB	Cross-attention	Vicuna	-	-	-	-
Emu	I+V+T→T	I/V: Eva-CLIP-IB	Cross-attention	LLaMA-13B	MLP	I: Stable Diffusion-1.5	-	-
DLP	I+T→T	I: CLIP/Eva-CLIP ViT	Q-Former+P-Former w/ Linear Projector	OPT/Flan-T5	-	-	-	-
BubbGPT	I+A+T→T+I <sub>h</sub>	I: CLIP/Eva-CLIP ViT; A: ImageBind	Q-Former w/ Linear Projector	Vicuna	-	-	-	-
ChatSpot	I+T→T	I: CLIP ViT-L/14	Linear Projector	Vicuna-7B/LLaMA	-	-	-	-
IDEFICS	I+T→T	I: OpenCLIP	Cross-attention	LLaMA	-	-	-	-
Qwen-VL-(Chat)	I+T→T	I: ViT@448 initialized from OpenCLIP's ViT-bigG	Cross-attention	Qwen-7B	-	-	1.4B <sup>†</sup>	50M <sup>†</sup>
LaViT	I+T→T	I: ViT	Cross-attention	LLaMA-7B	-	-	-	-
NExT-GPT	I+V+A+T→I+V+A+T	I/V/A: ImageBind	Linear Projector	Vicuna-7B	Tiny Transformer	I: Stable Diffusion; V: Zeroscope; A: AudioLDM	-	-
DreamLLM	I+T→T	I: CLIP ViT-L	Linear Projector	Vicuna	-	-	-	-
AnyMAL	I+V+A+T→T	I: CLIP ViT-L & ViT-G & DinoV2; V: InternVideo; A: CLAP	I/V: Cross-attention; A: Linear Projector	LLaMA-2	-	-	-	-
MiniGPT-5	I+T→T	I: Eva-CLIP ViT-G/14	Q-Former w/ Linear Projector	Vicuna-7B	Tiny Transformer w/ MLP	I: StableDiffusion-2	-	-
LLaVA-1.5	I+T→T	I: CLIP ViT-L@336	MLP	Vicuna-v1.5-7B/13B	-	-	0.6M	0.7M
MiniGPT-v2	I+T→T	I: Eva-CLIP ViT@448	Linear Projector	LLaMA-2-Chat-7B	-	-	-	-
CogVLM	I+T→T	I: Eva-2-CLIP ViT	MLP	Vicuna-v1.5-7B	-	-	-	-
Qwen-Audio	A+T→T	A: Whisper-L-v2	Linear Projector	Qwen-7B	-	-	-	-
DRESS	I+T→T	I: Eva-CLIP ViT-G/14	Linear Projector	Vicuna-v1.5-13B	-	-	-	-
X-InstructBLIP	I+V+A+3D+T→T	I/V: Eva-CLIP ViT-G/14; A: BEATs; 3D: ULIP-2	Q-Former w/ Linear Projector	Vicuna-v1.3-7B/13B	-	-	-	-
Code-2	I+V+A+T→I+V+A+T	I/V/A: ImageBind	MLP	LLaMA-2-Chat-7B	MLP	I: Stable Diffusion-2.1; V: Zeroscope-v2; A: AudioLDM-2	-	-
RLHF-V	I+T→T	I: BEIT-3	Linear Projector	Vicuna-v1-13B	-	-	-	-
Silkie	I+T→T	I: ViT initialized from OpenCLIP's ViT-bigG	Cross-attention	Qwen-7B	-	-	-	-
Lyrics	I+T→T	I: CLIP ViT-L/14 & Grounding-DINO-T w/ Swin-T & SAM-HQ w/ MAE & ViT & RAFT+ w/ Swin-B	MQ-Former w/ Linear Projection	Vicuna-13B	-	-	-	-
VILA	I+T→T	I: ViT@336	Linear Projector	LLaMA-2-7B/13B	-	-	50M	1M
IntenVL	I+V+T→T	I/V: IntenViT-6B; T: LLaMA-7B	Cross-attention w/ MLP	QLLaMA-8B & Vicuna-13B	-	-	-	-
ModaVerse	I+V+A+T→I+V+A+T	ImageBind	Linear Projector	LLaMA-2	MLP	I: Stable Diffusion; V: VideoFusion; A: AudioLDM	-	-
MM-Interleaved	I+T→I+T	I: CLIP ViT-L/14	Cross-attention	Vicuna-13B	Tiny Transformer	I: Stable Diffusion-2.1	-	-

Table 1: The summary of 43 mainstream MM-LLMs. I→O: Input to Output Modalities, I: Image, V: Video, A: Audio, 3D: Point Cloud, and T: Text. In Modality Encoder, “-L” represents Large, “-G” represents Giant, “/14” indicates a patch size of 14, and “@224” signifies an image resolution of  $224 \times 224$ . #\_PT and #\_IT represent the scale of dataset during MM PT and MM IT, respectively. <sup>†</sup> includes in-house data that is not publicly accessible.

conduct a comprehensive comparison of the architectures and training dataset scales for 43 of these SOTA MM-LLMs, as illustrated in Table 1. Next, we will summarize their developmental trends and briefly introduce the core contributions of some representative models in Appendix E.

**Trends in Existing MM-LLMs:** (1) Progressing from a dedicated emphasis on MM understanding to the generation of specific modalities and further evolving into any-to-any modality conversion (e.g., MiniGPT-4 → MiniGPT-5 → NExT-GPT); (2) Advancing from MM PT to SFT and then to RLHF, the training pipeline undergoes continuous refinement, striving to better align with human intent and enhance the model’s conversational interaction capabilities (e.g., BLIP-2 → InstructBLIP → DRESS); (3) Embracing Diversified Modal Extensions (e.g., BLIP-2 → X-LLM and InstructBLIP → X-InstructBLIP); (4) Incorporating a Higher-Quality Training Dataset (e.g., LLaVA → LLaVA-1.5); (5) Adopting a More Efficient Model Architecture, transitioning from complex Q- and P-Former input projector modules in BLIP-2 and DLP to a simpler yet effective linear projector in VILA.

## 5 Benchmarks and Performance

To offer a comprehensive performance comparison, we have compiled a table featuring major MM-LLMs across 18 VL benchmarks gathered

from various papers (Li et al., 2023e; Chen et al., 2023d,f; Lin et al., 2023), as shown in Table 2. The information of these benchmarks can be found in Appendix F. Next, we will extract essential training recipes that boost the effectiveness of MM-LLMs, drawing insights from SOTA models.

**Training Recipes** **Firstly**, higher image resolution can incorporate more visual details for the model, benefiting tasks that require fine-grained details. For example, LLaVA-1.5 and VILA employ a resolution of  $336 \times 336$ , while Qwen-VL and MiniGPT-v2 utilize  $448 \times 448$ . However, higher resolutions lead to longer token sequences, incurring additional training and inference costs. MiniGPT-v2 addresses this by concatenating 4 adjacent visual tokens in the embedding space to reduce length. Recently, Monkey (Li et al., 2023l) proposed a solution to enhance the resolution of input images without retraining a high-resolution visual encoder, utilizing only a low-resolution visual encoder, supporting resolutions up to  $1300 \times 800$ . To enhance the understanding of rich-text images, tables, and document content, DocPedia (Feng et al., 2023) introduced a method to increase the visual encoder resolution to  $2560 \times 2560$ , overcoming the limitations of poorly performing low resolutions in open-sourced ViT. **Secondly**, the incorporation of high-quality SFT data can significantly improve performance in specific tasks, as evidenced

Model	LLM Backbone	OKVQA	IconVQA	VQA <sup>v2</sup>	GQA	VizWiz	SQA <sup>1</sup>	VQA <sup>T</sup>	POPE	MME <sup>P</sup>	MME <sup>C</sup>	MMB	MMB <sup>CN</sup>	SEED <sup>1</sup>	LLaVA <sup>W</sup>	MM-Vet	QBench	HM	VSR
Flamingo	Chinchilla-7B	44.7	-	-	-	28.8	-	-	-	-	-	-	-	-	-	-	-	57.0	31.8
BLIP-2	Flan-T5xxl (13B)	45.9	40.6	65.0	44.7	19.6	61.0	42.5	85.3	1293.8	290.0	-	-	46.4	38.1	22.4	-	53.7	50.9
LLaVA	Vicuna-13B	54.4	43.0	-	41.3	-	-	38.9	-	-	-	-	-	-	-	-	-	-	51.2
MiniGPT-4	Vicuna-13B	37.5	37.6	-	30.8	-	-	19.4	-	-	-	-	-	-	-	-	-	-	41.6
InstructBLIP	Vicuna-7B	-	-	-	49.2	34.5	60.5	50.1	-	-	-	36.0	23.7	53.4	60.9	26.2	56.7	-	-
InstructBLIP	Vicuna-13B	-	44.8	-	49.5	33.4	63.1	50.7	78.9	1212.8	291.8	-	-	-	58.2	25.6	-	57.5	52.1
Shikra	Vicuna-13B	47.2	-	77.4 <sup>*</sup>	-	-	-	-	-	-	-	58.8	-	-	-	-	54.7	-	-
IDEFICS-9B	LLaMA-7B	-	-	50.9	38.4	35.5	-	25.9	-	-	-	48.2	25.2	-	-	-	-	-	-
IDEFICS-80B	LLaMA-65B	-	-	60.0	45.2	36.0	-	30.9	-	-	-	54.5	38.1	-	-	-	-	-	-
Qwen-VL	Qwen-7B	-	-	78.8 <sup>*</sup>	59.3 <sup>*</sup>	35.2	67.1	63.8	-	-	-	38.2	7.4	56.3	-	-	59.4	-	-
Qwen-VL-Chat	Qwen-7B	-	-	78.2 <sup>*</sup>	57.5 <sup>*</sup>	38.9	68.2	61.5	-	1487.5	360.7	60.6	56.7	58.2	-	-	-	-	-
LLaVA-1.5	Vicuna-1.5-7B	-	-	78.5 <sup>*</sup>	62.0 <sup>*</sup>	50.0	66.8	58.2	85.9	1510.7	316.1 <sup>†</sup>	64.3	58.3	58.6	63.4	30.5	58.7	-	-
+ShareGPT4V	Vicuna-1.5-7B	-	-	80.6	-	57.2	68.4	-	-	1567.4	376.4	68.8	62.2	69.7	72.6	37.6	63.4	-	-
LLaVA-1.5	Vicuna-1.5-13B	-	-	80.0 <sup>*</sup>	63.3 <sup>*</sup>	53.6	71.6	61.3	85.9	1531.3	295.4 <sup>†</sup>	67.7	63.6	61.6	70.7	35.4	62.1	-	-
MiniGPT-v2	LLaMA-2-Chat-7B	56.9	47.7	-	60.3	30.3	-	51.9	-	-	-	-	-	-	-	-	-	58.2	60.6
MiniGPT-v2-Chat	LLaMA-2-Chat-7B	55.9	49.4	-	58.8	42.4	-	52.3	-	-	-	-	-	-	-	-	-	59.5	63.3
VILA-7B	LLaMA-2-7B	-	-	79.9 <sup>*</sup>	62.3 <sup>*</sup>	57.8	68.2	64.4	85.5	1533.0	-	68.9	61.7	61.1	69.7	34.9	-	-	-
VILA-13B	LLaMA-2-13B	-	-	80.8 <sup>*</sup>	63.3 <sup>*</sup>	60.6	73.7	66.6	84.2	1570.1	-	70.3	64.3	62.8	73.0	38.8	-	-	-
+ShareGPT4V	LLaMA-2-13B	-	-	80.6 <sup>*</sup>	63.2 <sup>*</sup>	62.4	73.1	65.3	84.8	1556.5	-	70.8	65.4	61.4	78.4	45.7	-	-	-

Table 2: Comparison of mainstream MM-LLMs on 18 VL benchmarks. The **red** denotes the highest result, and the **blue** denotes the second highest result. <sup>†</sup> indicates ShareGPT4V’s (Chen et al., 2023f) re-implemented test results, which are missed in benchmarks or origin papers. \* indicates that training images are observed during training.

by the addition of ShareGPT4V data to LLaVA-1.5 and VILA-13B, as shown in Table 2. **Moreover**, VILA reveals several key findings: (1) Performing PEFT on the LLM Backbone promotes deep embedding alignment, crucial for ICL; (2) Interleaved Image-Text data proves beneficial, whereas Image-Text pairs alone are sub-optimal; (3) Re-blending text-only instruction data (*e.g.*, unnatural instruction (Honovich et al., 2022)) with image-text data during SFT not only addresses the degradation of text-only tasks but also enhances VL task accuracy.

## 6 Future Directions

In this section, we explore promising future directions for MM-LLMs across the following aspects:

**More Powerful Models** We can enhance the MM-LLMs’ strength from the following four key avenues: **(1) Expanding Modalities:** Current MM-LLMs mainly support the following modalities: image, video, audio, 3D, and text. However, the real world involves a broader range of modalities. Extending MM-LLMs to accommodate additional modalities (*e.g.*, web pages, heat maps, and figures&tables) will increase the model’s versatility, making it more universally applicable; **(2) Diversifying LLMs:** Incorporating various types and sizes of LLMs provides practitioners with the flexibility to select the most appropriate one based on their specific requirements; **(3) Improving MM IT Dataset Quality:** Current MM IT datasets have ample room for improvement and expansion. Diversifying the range of instructions can enhance the effectiveness of MM-LLMs in understanding and executing user commands. **(4) Strengthening MM Generation Capabilities:** Most current MM-LLMs are predominantly oriented towards MM understanding. Although some models have incor-

porated MM generation capabilities, the quality of generated responses may be constrained by the capacities of the LDMs. Exploring the integration of retrieval-based approaches (Asai et al., 2023; Gao et al., 2023a) holds significant promise in complementing the generative process, potentially enhancing the overall performance of the model.

**More Challenging Benchmarks** Existing benchmarks might not adequately challenge the capabilities of MM-LLMs, given that many datasets have previously appeared to varying degrees in the PT or IT sets. This implies that the models may have learned these tasks during training. Moreover, current benchmarks predominantly concentrate on the VL sub-field. Thus, it is crucial for the development of MM-LLMs to construct a more challenging, larger-scale benchmark that includes more modalities and uses a unified evaluation standard. For instance, GOAT-Bench (Lin et al., 2024b) is introduced to assess the capability of various MM-LLMs in discerning and responding to nuanced aspects of social abuse depicted in memes. MathVista (Lu et al., 2024) evaluates the math reasoning ability of MM-LLMs within visual contexts. Moreover, MMMU (Yue et al., 2023) and CMMMU (Zhang et al., 2024) have respectively introduced English and Chinese versions of the massive multi-discipline MM understanding and reasoning benchmark for expert artificial general intelligence. Fan et al. have also challenged MM-LLMs with multipanel VQA. BenchLMM (Cai et al., 2023) benchmarks the cross-style visual capability of MM-LLMs. Additionally, Liu et al. have conducted an in-depth study on the optical character recognition capabilities of MM-LLMs.

**Mobile/Lightweight Deployment** To deploy MM-LLMs on resource-constrained platforms and

495 achieve optimal performance meanwhile, such as  
496 low-power mobile and IoT devices, lightweight  
497 implementations are of paramount importance.  
498 A notable advancement in this realm is Mo-  
499 bileVLM (Chu et al., 2023a). This approach strate-  
500 gically downscales LLaMA, allowing for seam-  
501 less off-the-shelf deployment. MobileVLM fur-  
502 ther introduces a lightweight downsample pro-  
503 jector, consisting of fewer than 20 million pa-  
504 rameters, contributing to improved computational  
505 speed. Recently, there have been many simi-  
506 lar studies on lightweighting MM-LLMs, achiev-  
507 ing efficient computation and inference with com-  
508 parable performance or minimal loss, including  
509 TinyGPT-V (Yuan et al., 2023b), Vary-toy (Wei  
510 et al., 2024), Mobile-Agent (Wang et al., 2024b),  
511 MoE-LLaVA (Lin et al., 2024a), and MobileVLM  
512 V2 (Chu et al., 2024). Nevertheless, this avenue  
513 necessitates additional exploration for further ad-  
514 vancements in development.

515 **Embodied Intelligence** The embodied intelli-  
516 gence aims to replicate human-like perception and  
517 interaction with the surroundings by effectively  
518 understanding the environment, recognizing perti-  
519 nent objects, assessing their spatial relationships,  
520 and devising a comprehensive task plan (Firoozi  
521 et al., 2023). Embodied AI tasks, such as embod-  
522 ied planning, embodied visual question answer-  
523 ing, and embodied control, equip robots to au-  
524 tonomously implement extended plans by leverag-  
525 ing real-time observations. Some typical works in  
526 this area are PaLM-E (Driess et al., 2023) and Em-  
527 bodiedGPT (Mu et al., 2023). PaLM-E introduces  
528 a multi-embodiment agent through the training of  
529 a MM-LLM. Beyond functioning solely as an em-  
530 bodied decision maker, PaLM-E also demonstrates  
531 proficiency in handling general VL tasks. Em-  
532 bodiedGPT introduces an economically efficient  
533 method characterized by a CoT approach, enhanc-  
534 ing the capability of embodied agents to engage  
535 with the real world and establishing a closed loop  
536 that connects high-level planning with low-level  
537 control. While MM-LLM-based Embodied Intelli-  
538 gence has made advancements in integrating with  
539 robots, further exploration is needed to enhance the  
540 autonomy of robots.

541 **Continual Learning** Due to the large training  
542 costs associated with their massive scale, MM-  
543 LLMs are not amenable to frequent re-training.  
544 However, updates are necessary to endow MM-  
545 LLMs with new skills and keep them up-to-date

546 with rapidly evolving human knowledge (Wu et al.,  
547 2024). Thus, Continual Learning (CL) is needed to  
548 make the model flexible enough to efficiently and  
549 continually leverage emerging data while avoiding  
550 the substantial cost of retraining MM-LLMs. CL  
551 for MM-LLMs can be classified into two stages:  
552 continual PT and continual IT. Recently, a contin-  
553 ual MM IT benchmark has been proposed to con-  
554 tinuously fine-tune MM-LLMs for new MM tasks  
555 while maintaining superior performance on tasks  
556 learned during the original MM IT stage (He et al.,  
557 2023). It introduces two primary challenges: (1)  
558 catastrophic forgetting, where models forget previ-  
559 ous knowledge when learning new tasks (Robins,  
560 1995; McCloskey and Cohen, 1989; Goodfellow  
561 et al., 2013; Zhang et al., 2023d,c,b; Zheng et al.,  
562 2023a), and (2) negative forward transfer, indicat-  
563 ing that the performance of unseen tasks declines  
564 when learning new ones (Zheng et al., 2024)

565 **Mitigating Hallucination** Hallucinations entail  
566 generating textual descriptions of nonexistent ob-  
567 jects without visual cues, which manifest in diverse  
568 categories (Liu et al., 2024a) such as misjudgments  
569 and inaccuracies in descriptions. The origins of  
570 these hallucinations are multifaceted (Liu et al.,  
571 2024a), including biases and annotation errors in  
572 training data. Additionally, Skip  $\setminus n$  (Han et al.,  
573 2024) highlights semantic drift biases associated  
574 with paragraph separators, which can induce hal-  
575 lucinations when deliberately inserted. Current  
576 methods to mitigate these hallucinations involve  
577 leveraging self-feedback as visual cues (Lee et al.,  
578 2023). However, challenges persist, necessitat-  
579 ing nuanced discernment between accurate and hallu-  
580 cinatory outputs, as well as advancements in training  
581 methodologies to enhance output reliability.

## 582 7 Conclusion

583 In this paper, we have presented a comprehensive  
584 survey of MM-LLMs with a focus on recent ad-  
585 vancements. Initially, we categorize the model  
586 architecture into five components, providing a de-  
587 tailed overview of general design formulations and  
588 training pipelines. Subsequently, we introduce var-  
589 ious SOTA MM-LLMs, each distinguished by its  
590 specific formulations. Our survey also sheds light  
591 on their capabilities across diverse MM bench-  
592 marks and envisions future developments in this  
593 rapidly evolving field. We hope this survey can  
594 provide insights for researchers, contributing to the  
595 ongoing advancements in the MM-LLMs domain.



## 596 Limitations

597 In this paper, we embark on a comprehensive explo-  
598 ration of the current MM-LLMs landscape, present-  
599 ing a synthesis from diverse perspectives enriched  
600 by our insights. Acknowledging the dynamic na-  
601 ture of this field, it is plausible that certain aspects  
602 may have eluded our scrutiny, and recent advances  
603 might not be entirely encapsulated. To tackle this  
604 inherent challenge, we’ve established a dedicated  
605 website for real-time tracking, using crowdsourc-  
606 ing to capture the latest advancements. Our goal is  
607 for this platform to evolve into a continuous source  
608 of contributions propelling ongoing development  
609 in the field. Given the constraints of page limits,  
610 we are unable to delve into all technical details and  
611 have provided concise overviews of the core contri-  
612 butions of mainstream MM-LLMs. Looking ahead,  
613 we commit to vigilant monitoring and continual  
614 enhancement of relevant details on our website,  
615 incorporating fresh insights as they emerge.

## 616 References

617 2023. Bliva: A simple multimodal llm for better han-  
618 dling of text-rich visual questions. *arXiv preprint*  
619 *arXiv:2308.09936*.

620 Emanuele Aiello, Lili Yu, Yixin Nie, Armen Agha-  
621 janyan, and Barlas Oguz. 2023. Jointly Training  
622 Large Autoregressive Multimodal Models. *arXiv*  
623 *preprint arXiv:2309.15564*.

624 Hassan Akbari, Liangzhe Yuan, Rui Qian, Wei-Hong  
625 Chuang, Shih-Fu Chang, Yin Cui, and Boqing Gong.  
626 2021. Vatt: Transformers for multimodal self-  
627 supervised learning from raw video, audio and text.  
628 *Advances in Neural Information Processing Systems*,  
629 34:24206–24221.

630 Afra Feyza Akyürek, Ekin Akyürek, Aman Madaan,  
631 Ashwin Kalyan, Peter Clark, Derry Wijaya, and  
632 Niket Tandon. 2023. RL4F: Generating Natu-  
633 ral Language Feedback with Reinforcement Learn-  
634 ing for Repairing Model Outputs. *arXiv preprint*  
635 *arXiv:2305.08844*.

636 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,  
637 Antoine Miech, Iain Barr, Yana Hasson, Karel  
638 Lenc, Arthur Mensch, Katherine Millican, Malcolm  
639 Reynolds, et al. 2022. Flamingo: a visual language  
640 model for few-shot learning. *Advances in Neural*  
641 *Information Processing Systems*, 35:23716–23736.

642 Akari Asai, Sewon Min, Zexuan Zhong, and Danqi  
643 Chen. 2023. Retrieval-based language models and  
644 applications. In *Proceedings of the 61st Annual Meet-*  
645 *ing of the Association for Computational Linguistics*  
646 *(Volume 6: Tutorial Abstracts)*, pages 41–46.

Anas Awadalla, Irena Gao, Josh Gardner, Jack Hes-  
647 sel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe,  
648 Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al.  
649 2023. Openflamingo: An open-source framework for  
650 training large autoregressive vision-language models.  
651 *arXiv preprint arXiv:2308.01390*. 652

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,  
653 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei  
654 Huang, et al. 2023a. Qwen technical report. *arXiv*  
655 *preprint arXiv:2309.16609*. 656

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang,  
657 Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou,  
658 and Jingren Zhou. 2023b. Qwen-VL: A Frontier  
659 Large Vision-Language Model with Versatile Abili-  
660 ties. *CoRR*, abs/2308.12966. 661

Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zis-  
662 serman. 2021. Frozen in time: A joint video and  
663 image encoder for end-to-end retrieval. In *Proceed-*  
664 *ings of the IEEE/CVF International Conference on*  
665 *Computer Vision*, pages 1728–1738. 666

Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. 2022. 667  
668 Analytic-DPM: an Analytic Estimate of the Optimal  
669 Reverse Variance in Diffusion Probabilistic Models.  
670 In *International Conference on Learning Representa-*  
671 *tions*. 671

Rohan Bavishi, Erich Elsen, Curtis Hawthorne,  
672 Maxwell Nye, Augustus Odena, Arushi Somani, and  
673 Sağnak Taşırlar. 2023. [Introducing our Multimodal](#)  
674 [Models](#). 675

Ali Furkan Biten, Ron Litman, Yusheng Xie, Srikar  
676 Appalaraju, and R Manmatha. 2022. Latr: Layout-  
677 aware transformer for scene-text vqa. In *Proceedings*  
678 *of the IEEE/CVF conference on computer vision and*  
679 *pattern recognition*, pages 16548–16558. 680

Andy Brock, Soham De, Samuel L Smith, and Karen Si-  
681 monyan. 2021. High-performance large-scale image  
682 recognition without normalization. In *International*  
683 *Conference on Machine Learning*, pages 1059–1071.  
684 PMLR. 685

Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
686 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
687 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
688 Askell, et al. 2020. Language models are few-shot  
689 learners. *Advances in neural information processing*  
690 *systems*, 33:1877–1901. 691

Minwoo Byeon, Beomhee Park, Haecheon Kim,  
692 Sungjun Lee, Woonhyuk Baek, and Saehoon Kim.  
693 2022. [Coyo-700m: Image-text pair dataset](#). 694

Fabian Caba Heilbron, Victor Escorcia, Bernard  
695 Ghanem, and Juan Carlos Nieves. 2015. Activitynet:  
696 A large-scale video benchmark for human activity  
697 understanding. In *Proceedings of the ieee conference*  
698 *on computer vision and pattern recognition*, pages  
699 961–970. 700



811	Xiangxiang Chu, Limeng Qiao, Xinyu Zhang, Shuang Xu, Fei Wei, Yang Yang, Xiaofei Sun, Yiming Hu, Xinyang Lin, Bo Zhang, et al. 2024. MobileVLM V2: Faster and Stronger Baseline for Vision Language Model. <i>arXiv preprint arXiv:2402.03766</i> .	Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In <i>International Conference on Learning Representations</i> .	867
812			868
813			869
814			870
815			
816	Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. 2023b. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. <i>arXiv preprint arXiv:2311.07919</i> .	Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. 2023. Palm-e: An embodied multimodal language model. <i>arXiv preprint arXiv:2303.03378</i> .	871
817			872
818			873
819			874
820			875
821	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. <i>arXiv preprint arXiv:2210.11416</i> .	Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. 2022a. A Survey of Vision-Language Pre-Trained Models. In <i>Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022</i> , pages 5436–5443.	876
822			877
823			878
824			879
825			880
826	XTuner Contributors. 2023. XTuner: A Toolkit for Efficiently Fine-tuning LLM. <a href="https://github.com/InternLM/xtuner">https://github.com/InternLM/xtuner</a> .	Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022b. GLM: General Language Model Pretraining with Autoregressive Blank Infilling. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 320–335.	881
827			882
828			883
829	Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zichong Yang, Kuei-Da Liao, et al. 2024. A survey on multimodal large language models for autonomous driving. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pages 958–979.	Erich Elsen, Augustus Odena, Maxwell Nye, Sağnak Taşırlar, Tri Dao, Curtis Hawthorne, Deepak Moparthy, and Arushi Somani. 2023. <a href="#">Releasing Persimmon-8B</a> .	884
830			885
831			886
832			887
833			888
834			889
835			890
836	Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven C. H. Hoi. 2023. InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	Yue Fan, Jing Gu, Kaiwen Zhou, Qianqi Yan, Shan Jiang, Ching-Chen Kuo, Xinze Guan, and Xin Eric Wang. 2024. Muffin or Chihuahua? Challenging Large Vision-Language Models with Multipanel VQA. <i>arXiv preprint arXiv:2401.15847</i> .	891
837			892
838			893
839			894
840			895
841			896
842			897
843	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. <i>arXiv preprint arXiv:2305.14314</i> .	Han Fang, Pengfei Xiong, Luhui Xu, and Yu Chen. 2021. Clip2video: Mastering video-text retrieval via image clip. <i>arXiv preprint arXiv:2106.11097</i> .	898
844			899
845			900
846	Lin hao Dong and Bo Xu. 2020. Cif: Continuous integrate-and-fire for end-to-end speech recognition. In <i>ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 6079–6083. IEEE.	Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. 2023. Eva: Exploring the limits of masked visual representation learning at scale. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 19358–19369.	901
847			902
848			903
849			904
850			905
851	Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu Zhou, Haoran Wei, et al. 2024a. Dreamllm: Synergistic multimodal comprehension and creation. In <i>The Twelfth International Conference on Learning Representations</i> .	Hao Feng, Qi Liu, Hao Liu, Wengang Zhou, Houqiang Li, and Can Huang. 2023. DocPedia: Unleashing the Power of Large Multimodal Model in the Frequency Domain for Versatile Document Understanding. <i>arXiv preprint arXiv:2311.11810</i> .	906
852			907
853			908
854			909
855			910
856			911
857	Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang Zhang, Haodong Duan, Maosong Cao, et al. 2024b. InternLM-XComposer2: Mastering Free-form Text-Image Composition and Comprehension in Vision-Language Large Model. <i>arXiv preprint arXiv:2401.16420</i> .	Roya Firoozi, Johnathan Tucker, Stephen Tian, Anirudha Majumdar, Jiankai Sun, Weiyu Liu, Yuke Zhu, Shuran Song, Ashish Kapoor, Karol Hausman, et al. 2023. Foundation Models in Robotics: Applications, Challenges, and the Future. <i>arXiv preprint arXiv:2312.07843</i> .	912
858			913
859			914
860			915
861			916
862			917
863			918
864	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias	Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. 2023. Mme: A comprehensive evaluation benchmark for multimodal large language models. <i>arXiv preprint arXiv:2306.13394</i> .	919
865			920
866			921
			922
			923

924	Chin-Lun Fu, Zih-Ching Chen, Yun-Ru Lee, and Hung-Yi Lee. 2022. AdapterBias: Parameter-efficient Token-dependent Representation Shift for Adapters in NLP Tasks. In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 2608–2621.	<i>IEEE conference on computer vision and pattern recognition</i> , pages 6904–6913.	980 981
930	Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. 2023. Datacomp: In search of the next generation of multimodal datasets. <i>arXiv preprint arXiv:2304.14108</i> .	Jiayi Gu, Xiaojun Meng, Guansong Lu, Lu Hou, Niu Minzhe, Xiaodan Liang, Lewei Yao, Runhui Huang, Wei Zhang, Xin Jiang, et al. 2022. Wukong: A 100 million large-scale chinese cross-modal pre-training benchmark. <i>Advances in Neural Information Processing Systems</i> , 35:26418–26431.	982 983 984 985 986 987
936	Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2023. Bias and fairness in large language models: A survey. <i>arXiv preprint arXiv:2309.00770</i> .	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. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 3608–3617.	988 989 990 991 992 993
941	Peng Gao, Renrui Zhang, Chris Liu, Longtian Qiu, Siyuan Huang, Weifeng Lin, Shitian Zhao, Shijie Geng, Ziyi Lin, Peng Jin, et al. 2024. SPHINX-X: Scaling Data and Parameters for a Family of Multi-modal Large Language Models. <i>arXiv preprint arXiv:2402.05935</i> .	Minglun Han, Feilong Chen, Jing Shi, Shuang Xu, and Bo Xu. 2023. Knowledge Transfer from Pre-trained Language Models to Cif-based Speech Recognizers via Hierarchical Distillation. <i>arXiv preprint arXiv:2301.13003</i> .	994 995 996 997 998
947	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023a. Retrieval-augmented generation for large language models: A survey. <i>arXiv preprint arXiv:2312.10997</i> .	Minglun Han, Linhao Dong, Zhenlin Liang, Meng Cai, Shiyu Zhou, Zejun Ma, and Bo Xu. 2022. Improving end-to-end contextual speech recognition with fine-grained contextual knowledge selection. In <i>ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 8532–8536. IEEE.	999 1000 1001 1002 1003 1004 1005
952	Zhi Gao, Yuntao Du, Xintong Zhang, Xiaojian Ma, Wenjuan Han, Song-Chun Zhu, and Qing Li. 2023b. CLOVA: A Closed-Loop Visual Assistant with Tool Usage and Update. <i>arXiv preprint arXiv:2312.10908</i> .	Zongbo Han, Zechen Bai, Haiyang Mei, Qianli Xu, Changqing Zhang, and Mike Zheng Shou. 2024. Skip $\setminus n$ : A simple method to reduce hallucination in large vision-language models. <i>arXiv preprint arXiv:2402.01345</i> .	1006 1007 1008 1009 1010
957	Yuying Ge, Yixiao Ge, Ziyun Zeng, Xintao Wang, and Ying Shan. 2023. Planting a seed of vision in large language model. <i>arXiv preprint arXiv:2307.08041</i> .	Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, and Dong Yu. 2024. WebVoyager: Building an End-to-End Web Agent with Large Multimodal Models. <i>arXiv preprint arXiv:2401.13919</i> .	1011 1012 1013 1014 1015
960	Rohit Girdhar, Alaeldin El-Nouby, Zhuang Liu, Man- nat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. 2023. Imagebind: One embed- ding space to bind them all. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pat- tern Recognition</i> , pages 15180–15190.	Jinghan He, Haiyun Guo, Ming Tang, and Jinqiao Wang. 2023. Continual instruction tuning for large multi-modal models. <i>arXiv preprint arXiv:2311.16206</i> .	1016 1017 1018
966	Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. 2023. Multimodal-gpt: A vision and language model for dialogue with humans. <i>arXiv preprint arXiv:2305.04790</i> .	Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg- Kirkpatrick, and Graham Neubig. 2021. Towards a Unified View of Parameter-Efficient Transfer Learn- ing. In <i>International Conference on Learning Repre- sentations</i> .	1019 1020 1021 1022 1023
971	Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. 2013. An em- pirical investigation of catastrophic forgetting in gradient-based neural networks. <i>arXiv preprint arXiv:1312.6211</i> .	Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Pi- otr Dollár, and Ross Girshick. 2022. Masked autoen- coders are scalable vision learners. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 16000–16009.	1024 1025 1026 1027 1028
976	Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In <i>Proceedings of the</i>	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recog- nition. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 770– 778.	1029 1030 1031 1032 1033



1034	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. Training compute-optimal large language models. <i>arXiv preprint arXiv:2203.15556</i> .	1090
1035		1091
1036		
1037		
1038		
1039		
1040	Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. 2023. Cogagent: A visual language model for gui agents. <i>arXiv preprint arXiv:2312.08914</i> .	
1041		
1042		
1043		
1044		
1045	Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. 2022. Unnatural instructions: Tuning language models with (almost) no human labor. <i>arXiv preprint arXiv:2212.09689</i> .	
1046		
1047		
1048		
1049	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. In <i>International Conference on Machine Learning</i> , pages 2790–2799. PMLR.	
1050		
1051		
1052		
1053		
1054		
1055	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. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 29:3451–3460.	
1056		
1057		
1058		
1059		
1060		
1061	Anwen Hu, Yaya Shi, Haiyang Xu, Jiabo Ye, Qinghao Ye, Ming Yan, Chenliang Li, Qi Qian, Ji Zhang, and Fei Huang. 2023a. mPLUG-PaperOwl: Scientific Diagram Analysis with the Multimodal Large Language Model. <i>arXiv preprint arXiv:2311.18248</i> .	
1062		
1063		
1064		
1065		
1066	Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2021. LoRA: Low-Rank Adaptation of Large Language Models. In <i>International Conference on Learning Representations</i> .	
1067		
1068		
1069		
1070		
1071	Jinyi Hu, Yuan Yao, Chongyi Wang, Shan Wang, Yinxu Pan, Qianyu Chen, Tianyu Yu, Hanghao Wu, Yue Zhao, Haoye Zhang, et al. 2023b. Large multilingual models pivot zero-shot multimodal learning across languages. <i>arXiv preprint arXiv:2308.12038</i> .	
1072		
1073		
1074		
1075		
1076	Jiaxing Huang, Jingyi Zhang, Kai Jiang, Han Qiu, and Shijian Lu. 2023a. Visual Instruction Tuning towards General-Purpose Multimodal Model: A Survey. <i>arXiv preprint arXiv:2312.16602</i> .	
1077		
1078		
1079		
1080	Rongjie Huang, Mingze Li, Dongchao Yang, Jia-tong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, et al. 2023b. Audiogpt: Understanding and generating speech, music, sound, and talking head. <i>arXiv preprint arXiv:2304.12995</i> .	
1081		
1082		
1083		
1084		
1085		
1086	Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, et al. 2023c. Language is not all you need: Aligning	
1087		
1088		
1089		
	perception with language models. <i>arXiv preprint arXiv:2302.14045</i> .	1092
		1093
		1094
		1095
		1096
	Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 6700–6709.	1097
		1098
	IDEFICS. 2023. <a href="#">Introducing IDEFICS: An Open Reproduction of State-of-the-Art Visual Language Model</a> .	1099
		1100
		1101
		1102
		1103
		1104
	Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Daniel Simig, Ping Yu, Kurt Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, et al. 2022. Opt-impl: Scaling language model instruction meta learning through the lens of generalization. <i>arXiv preprint arXiv:2212.12017</i> .	1105
		1106
		1107
		1108
	Jitesh Jain, Jianwei Yang, and Humphrey Shi. 2023. Vcoder: Versatile vision encoders for multimodal large language models. <i>arXiv preprint arXiv:2312.14233</i> .	1109
		1110
		1111
		1112
		1113
		1114
	Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In <i>International conference on machine learning</i> , pages 4904–4916. PMLR.	1115
		1116
		1117
		1118
		1119
	Yiren Jian, Chongyang Gao, and Soroush Vosoughi. 2023. Bootstrapping Vision-Language Learning with Decoupled Language Pre-training. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1120
		1121
		1122
		1123
		1124
		1125
	Yang Jin, Kun Xu, Liwei Chen, Chao Liao, Jianchao Tan, Bin Chen, Chenyi Lei, An Liu, Chengru Song, Xiaoqiang Lei, et al. 2024. Unified language-vision pretraining with dynamic discrete visual tokenization. In <i>The Twelfth International Conference on Learning Representations</i> .	1126
		1127
		1128
		1129
		1130
	Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. 2018. Dvqa: Understanding data visualizations via question answering. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 5648–5656.	1131
		1132
		1133
		1134
	Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. 2021. Compacter: Efficient low-rank hypercomplex adapter layers. <i>Advances in Neural Information Processing Systems</i> , 34:1022–1035.	1135
		1136
		1137
		1138
		1139
		1140
	Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. 2014. Referitgame: Referring to objects in photographs of natural scenes. In <i>Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)</i> , pages 787–798.	1141
		1142
		1143
		1144
	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In <i>Proceedings of NAACL-HLT</i> , pages 4171–4186.	

1145	Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. 2020. The hateful memes challenge: Detecting hate speech in multimodal memes. <i>Advances in neural information processing systems</i> , 33:2611–2624.	1199
1146		1200
1147		1201
1148		1202
1149		
1150		
1151	Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. <i>arXiv preprint arXiv:1312.6114</i> .	
1152		
1153		
1154	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything. <i>arXiv preprint arXiv:2304.02643</i> .	
1155		
1156		
1157		
1158		
1159	Jing Yu Koh, Daniel Fried, and Ruslan Salakhutdinov. 2023a. Generating images with multimodal language models. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	
1160		
1161		
1162		
1163	Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. 2023b. Grounding language models to images for multimodal inputs and outputs. In <i>International Conference on Machine Learning</i> , pages 17283–17300. PMLR.	
1164		
1165		
1166		
1167		
1168	Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. <i>International journal of computer vision</i> , 123:32–73.	
1169		
1170		
1171		
1172		
1173		
1174	Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. 2023. Lisa: Reasoning segmentation via large language model. <i>arXiv preprint arXiv:2308.00692</i> .	
1175		
1176		
1177		
1178	Hugo Laurençon, Lucile Saulnier, Leo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M Rush, Douwe Kiela, et al. 2023. OBELICS: An Open Web-Scale Filtered Dataset of Interleaved Image-Text Documents. In <i>Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	
1179		
1180		
1181		
1182		
1183		
1184		
1185		
1186	Seongyun Lee, Sue Hyun Park, Yongrae Jo, and Minjoon Seo. 2023. Volcano: mitigating multimodal hallucination through self-feedback guided revision. <i>arXiv preprint arXiv:2311.07362</i> .	
1187		
1188		
1189		
1190	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The Power of Scale for Parameter-Efficient Prompt Tuning. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 3045–3059.	
1191		
1192		
1193		
1194		
1195	Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. 2023a. Mimic-it: Multi-modal in-context instruction tuning. <i>arXiv preprint arXiv:2306.05425</i> .	
1196		
1197		
1198		
	Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023b. Otter: A multi-modal model with in-context instruction tuning. <i>arXiv preprint arXiv:2305.03726</i> .	1203
		1204
		1205
		1206
	Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023c. Seed-bench: Benchmarking multimodal llms with generative comprehension. <i>arXiv preprint arXiv:2307.16125</i> .	1207
		1208
		1209
		1210
		1211
		1212
	Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023d. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. <i>arXiv preprint arXiv:2306.00890</i> .	1213
		1214
		1215
		1216
		1217
		1218
	Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. 2023e. BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , pages 19730–19742.	1219
		1220
		1221
		1222
		1223
	Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In <i>International Conference on Machine Learning</i> , pages 12888–12900. PMLR.	1224
		1225
		1226
		1227
		1228
		1229
	Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. 2021. Align before fuse: Vision and language representation learning with momentum distillation. <i>Advances in neural information processing systems</i> , 34:9694–9705.	1230
		1231
		1232
		1233
	KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023f. Videochat: Chat-centric video understanding. <i>arXiv preprint arXiv:2305.06355</i> .	1234
		1235
		1236
		1237
		1238
	Lei Li, Zhihui Xie, Mukai Li, Shunian Chen, Peiyi Wang, Liang Chen, Yazheng Yang, Benyou Wang, and Lingpeng Kong. 2023g. Silkie: Preference Distillation for Large Visual Language Models. <i>arXiv preprint arXiv:2312.10665</i> .	1239
		1240
		1241
		1242
		1243
	Lei Li, Yuwei Yin, Shicheng Li, Liang Chen, Peiyi Wang, Shuhuai Ren, Mukai Li, Yazheng Yang, Jingjing Xu, Xu Sun, et al. 2023h. M <sup>3</sup> IT: A Large-Scale Dataset towards Multi-Modal Multilingual Instruction Tuning. <i>arXiv preprint arXiv:2306.04387</i> .	1244
		1245
		1246
	Mukai Li, Lei Li, Yuwei Yin, Masood Ahmed, Zhen-guang Liu, and Qi Liu. 2024a. Red teaming visual language models. <i>arXiv preprint arXiv:2401.12915</i> .	1247
		1248
		1249
		1250
	Xiang Lisa Li and Percy Liang. 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4582–4597.	1251
		1252
		1253

1254	Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang,	Fangyu Liu, Guy Emerson, and Nigel Collier. 2023a.	1309
1255	Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong	Visual spatial reasoning. <i>Transactions of the Associ-</i>	1310
1256	Hu, Li Dong, Furu Wei, et al. 2020. Oscar: Object-	<i>ation for Computational Linguistics</i> , 11:635–651.	1311
1257	semantics aligned pre-training for vision-language		
1258	tasks. In <i>Computer Vision–ECCV 2020: 16th Euro-</i>	Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen,	1312
1259	<i>pean Conference, Glasgow, UK, August 23–28, 2020,</i>	Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li,	1313
1260	<i>Proceedings, Part XXX 16</i> , pages 121–137. Springer.	and Wei Peng. 2024a. A survey on hallucination	1314
		in large vision-language models. <i>arXiv preprint</i>	1315
1261	Yanda Li, Chi Zhang, Gang Yu, Zhibin Wang, Bin	<i>arXiv:2402.00253</i> .	1316
1262	Fu, Guosheng Lin, Chunhua Shen, Ling Chen, and		
1263	Yunchao Wei. 2023i. Stablelava: Enhanced visual	Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo	1317
1264	instruction tuning with synthesized image-dialogue	Liu, Danilo P. Mandic, Wenwu Wang, and Mark D.	1318
1265	data. <i>arXiv preprint arXiv:2308.10253</i> .	Plumbly. 2023b. AudioLDM: Text-to-Audio Gener-	1319
		ation with Latent Diffusion Models. In <i>International</i>	1320
1266	Yanwei Li, Chengyao Wang, and Jiaya Jia. 2023j.	<i>Conference on Machine Learning, ICML 2023, 23-</i>	1321
1267	LLaMA-VID: An Image is Worth 2 Tokens in Large	<i>29 July 2023, Honolulu, Hawaii, USA</i> , pages 21450–	1322
1268	Language Models. <i>arXiv preprint arXiv:2311.17043</i> .	21474.	1323
1269	Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang,	Haohe Liu, Qiao Tian, Yi Yuan, Xubo Liu, Xinhao	1324
1270	Wayne Xin Zhao, and Ji-Rong Wen. 2023k. Eval-	Mei, Qiuqiang Kong, Yuping Wang, Wenwu Wang,	1325
1271	uating object hallucination in large vision-language	Yuxuan Wang, and Mark D. Plumbly. 2023c. Audi-	1326
1272	models. <i>arXiv preprint arXiv:2305.10355</i> .	oLDM 2: Learning Holistic Audio Generation with	1327
		Self-supervised Pretraining. <i>CoRR</i> , abs/2308.05734.	1328
1273	Zeju Li, Chao Zhang, Xiaoyan Wang, Ruilong Ren, Yi-		
1274	fan Xu, Ruifei Ma, and Xiangde Liu. 2024b. 3DMIT:	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae	1329
1275	3D Multi-modal Instruction Tuning for Scene Under-	Lee. 2023d. Improved Baselines with Visual Instruc-	1330
1276	standing. <i>arXiv preprint arXiv:2401.03201</i> .	tion Tuning. In <i>NeurIPS 2023 Workshop on Instruc-</i>	1331
		<i>tion Tuning and Instruction Following</i> .	1332
1277	Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo		
1278	Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuan-	1333
1279	Xiang Bai. 2023l. Monkey: Image Resolution and	han Zhang, Sheng Shen, and Yong Jae Lee. 2024b.	1334
1280	Text Label Are Important Things for Large Multi-	LLaVA-NeXT: Improved reasoning, OCR, and world	1335
1281	modal Models. <i>arXiv preprint arXiv:2311.06607</i> .	knowledge.	1336
1282	Zhaowei Li, Qi Xu, Dong Zhang, Hang Song, Yiqing	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	1337
1283	Cai, Qi Qi, Ran Zhou, Junting Pan, Zefeng Li,	Lee. 2023e. Visual Instruction Tuning. In <i>Thirty-</i>	1338
1284	Van Tu Vu, et al. 2024c. LEGO: Language Enhanced	<i>seventh Conference on Neural Information Process-</i>	1339
1285	Multi-modal Grounding Model. <i>arXiv preprint</i>	<i>ing Systems</i> .	1340
1286	<i>arXiv:2401.06071</i> .		
1287	Bin Lin, Zhenyu Tang, Yang Ye, Jiayi Cui, Bin	Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng	1341
1288	Zhu, Peng Jin, Junwu Zhang, Munan Ning, and	Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su,	1342
1289	Li Yuan. 2024a. MoE-LLaVA: Mixture of Experts	Jun Zhu, et al. 2023f. Llava-plus: Learning to use	1343
1290	for Large Vision-Language Models. <i>arXiv preprint</i>	tools for creating multimodal agents. <i>arXiv preprint</i>	1344
1291	<i>arXiv:2401.15947</i> .	<i>arXiv:2311.05437</i> .	1345
1292	Hongzhan Lin, Ziyang Luo, Bo Wang, Ruichao Yang,	Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengx-	1346
1293	and Jing Ma. 2024b. GOAT-Bench: Safety Insights	iao Du, Zhilin Yang, and Jie Tang. 2022. P-tuning:	1347
1294	to Large Multimodal Models through Meme-Based	Prompt tuning can be comparable to fine-tuning	1348
1295	Social Abuse. <i>arXiv preprint arXiv:2401.01523</i> .	across scales and tasks. In <i>Proceedings of the 60th</i>	1349
		<i>Annual Meeting of the Association for Computational</i>	1350
1296	Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo	<i>Linguistics (Volume 2: Short Papers)</i> , pages 61–68.	1351
1297	Molchanov, Andrew Tao, Huizi Mao, Jan Kautz,		
1298	Mohammad Shoeybi, and Song Han. 2023. VILA:	Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam,	1352
1299	On Pre-training for Visual Language Models. <i>arXiv</i>	Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021a.	1353
1300	<i>preprint arXiv:2312.07533</i> .	P-tuning v2: Prompt tuning can be comparable to	1354
		fine-tuning universally across scales and tasks. <i>arXiv</i>	1355
1301	Tsung-Yi Lin, Michael Maire, Serge Belongie, James	<i>preprint arXiv:2110.07602</i> .	1356
1302	Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,		
1303	and C Lawrence Zitnick. 2014. Microsoft coco:	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li,	1357
1304	Common objects in context. In <i>Computer Vision–</i>	Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi	1358
1305	<i>ECCV 2014: 13th European Conference, Zurich,</i>	Wang, Conghui He, Ziwei Liu, et al. 2023g. Mm-	1359
1306	<i>Switzerland, September 6-12, 2014, Proceedings,</i>	<i>bench: Is your multi-modal model an all-around</i>	1360
1307	<i>Part V 13</i> , pages 740–755. Springer.	<i>player? arXiv preprint arXiv:2307.06281</i> .	1361
1308	LinkSoul-AI. 2023. <a href="#">Chinese-LLaVA</a> .		





1476	Ziyi Ni, Minglun Han, Feilong Chen, Linghui Meng, Jing Shi, Pin Lv, and Bo Xu. 2024. VILAS: Exploring the Effects of Vision and Language Context in Automatic Speech Recognition. In <i>ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> . IEEE.	1531
1477		1532
1478		1533
1479		1534
1480		1535
1481		1536
1482	OpenAI. 2022. <i>OpenAI: Introducing ChatGPT</i> .	1537
1483	OpenAI. 2023. <i>GPT-4 Technical Report</i> .	1538
1484		1539
1485	Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. 2023. Dinov2: Learning robust visual features without supervision. <i>arXiv preprint arXiv:2304.07193</i> .	1540
1486		1541
1487		1542
1488		1543
1489		1544
1490	Vicente Ordonez, Girish Kulkarni, and Tamara Berg. 2011. Im2text: Describing images using 1 million captioned photographs. <i>Advances in neural information processing systems</i> , 24.	1545
1491		1546
1492		1547
1493		1548
1494		1549
1495	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	1550
1496		1551
1497		1552
1498		1553
1499		1554
1500	Xichen Pan, Li Dong, Shaohan Huang, Zhiliang Peng, Wenhui Chen, and Furu Wei. 2023. Kosmos-g: Generating images in context with multimodal large language models. <i>arXiv preprint arXiv:2310.02992</i> .	1555
1501		1556
1502		1557
1503		1558
1504	Artemis Panagopoulou, Le Xue, Ning Yu, Junnan Li, Dongxu Li, Shafiq Joty, Ran Xu, Silvio Savarese, Caiming Xiong, and Juan Carlos Niebles. 2023. X-InstructBLIP: A Framework for aligning X-Modal instruction-aware representations to LLMs and Emergent Cross-modal Reasoning. <i>arXiv preprint arXiv:2311.18799</i> .	1559
1505		1560
1506		1561
1507		1562
1508		1563
1509		1564
1510		1565
1511	Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding Multimodal Large Language Models to the World. <i>arXiv preprint arXiv:2306.14824</i> .	1566
1512		1567
1513		1568
1514		1569
1515		1570
1516		1571
1517	Renjie Pi, Jiahui Gao, Shizhe Diao, Rui Pan, Hanze Dong, Jipeng Zhang, Lewei Yao, Jianhua Han, Hang Xu, and Lingpeng Kong Tong Zhang. 2023. DetGPT: Detect What You Need via Reasoning. <i>arXiv preprint arXiv:2305.14167</i> .	1572
1518		1573
1519		1574
1520		1575
1521		1576
1522		1577
1523		1578
1524		1579
1525		1580
1526		1581
1527		1582
1528		1583
1529		1584
1530		1585
		1586
	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastri, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.	
	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust Speech Recognition via Large-Scale Weak Supervision. In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , pages 28492–28518.	
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>The Journal of Machine Learning Research</i> , 21(1):5485–5551.	
	Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Erix Xing, Ming-Hsuan Yang, and Fahad S Khan. 2023. Glamm: Pixel grounding large multimodal model. <i>arXiv preprint arXiv:2311.03356</i> .	
	Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. 2017. Learning multiple visual domains with residual adapters. <i>Advances in neural information processing systems</i> , 30.	
	Zhongwei Ren, Zhicheng Huang, Yunchao Wei, Yao Zhao, Dongmei Fu, Jiashi Feng, and Xiaojie Jin. 2023. PixelLM: Pixel Reasoning with Large Multimodal Model. <i>arXiv preprint arXiv:2312.02228</i> .	
	Anthony Robins. 1995. Catastrophic forgetting, rehearsal and pseudorehearsal. <i>Connection Science</i> , 7(2):123–146.	
	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 10684–10695.	
	Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In <i>Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18</i> , pages 234–241. Springer.	
	Ludan Ruan and Qin Jin. 2022. Survey: Transformer based video-language pre-training. <i>AI Open</i> , 3:1–13.	
	Paul K Rubenstein, Chulayuth Asawaroengchai, Duc Dung Nguyen, Ankur Bapna, Zalán Borsos, Félix de Chaumont Quitry, Peter Chen, Dalia El Badawy, Wei Han, Eugene Kharitonov, et al. 2023. AudioPaLM: A Large Language Model That Can Speak and Listen. <i>arXiv preprint arXiv:2306.12925</i> .	

1587	Salesforce. 2022. <a href="#">Ulip</a> .	1640
1588	Christoph Schuhmann, Romain Beaumont, Richard	1641
1589	Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti,	1642
1590	Theo Coombes, Aarush Katta, Clayton Mullis,	1643
1591	Mitchell Wortsman, et al. 2022. Laion-5b: An open	1644
1592	large-scale dataset for training next generation image-	
1593	text models. <i>Advances in Neural Information Pro-</i>	
1594	<i>cessing Systems</i> , 35:25278–25294.	
1595	Christoph Schuhmann, Andreas Köpf, Richard Vencu,	
1596	Theo Coombes, and Romain Beaumont. 2022b.	
1597	<a href="#">Laion coco: 600m synthetic captions from laion2b-</a>	
1598	<a href="#">en</a> .	
1599	Christoph Schuhmann, Richard Vencu, Romain Beau-	
1600	mont, Robert Kaczmarczyk, Clayton Mullis, Aarush	
1601	Katta, Theo Coombes, Jenia Jitsev, and Aran Komat-	
1602	suzaki. 2021. Laion-400m: Open dataset of clip-	
1603	filtered 400 million image-text pairs. <i>arXiv preprint</i>	
1604	<i>arXiv:2111.02114</i> .	
1605	Dustin Schwenk, Apoorv Khandelwal, Christopher	
1606	Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022.	
1607	A-okvqa: A benchmark for visual question answer-	
1608	ing using world knowledge. In <i>European Conference</i>	
1609	<i>on Computer Vision</i> , pages 146–162. Springer.	
1610	Piyush Sharma, Nan Ding, Sebastian Goodman, and	
1611	Radu Soricut. 2018. Conceptual captions: A cleaned,	
1612	hypernymed, image alt-text dataset for automatic im-	
1613	age captioning. In <i>Proceedings of the 56th Annual</i>	
1614	<i>Meeting of the Association for Computational Lin-</i>	
1615	<i>guistics (Volume 1: Long Papers)</i> , pages 2556–2565.	
1616	Weizhou Shen, Chenliang Li, Hongzhan Chen, Ming	
1617	Yan, Xiaojun Quan, Hehong Chen, Ji Zhang, and Fei	
1618	Huang. 2024. Small llms are weak tool learners: A	
1619	multi-llm agent. <i>arXiv preprint arXiv:2401.07324</i> .	
1620	Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li,	
1621	Weiming Lu, and Yueting Zhuang. 2023. Hugging-	
1622	gpt: Solving ai tasks with chatgpt and its friends in	
1623	huggingface. <i>arXiv preprint arXiv:2303.17580</i> .	
1624	Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and	
1625	Amanpreet Singh. 2020. Textcaps: a dataset for im-	
1626	age captioning with reading comprehension. In <i>Com-</i>	
1627	<i>puter Vision—ECCV 2020: 16th European Confer-</i>	
1628	<i>ence, Glasgow, UK, August 23–28, 2020, Proceed-</i>	
1629	<i>ings, Part II 16</i> , pages 742–758. Springer.	
1630	Amanpreet Singh, Vivek Natarajan, Meet Shah,	
1631	Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh,	
1632	and Marcus Rohrbach. 2019. Towards vqa models	
1633	that can read. In <i>Proceedings of the IEEE/CVF con-</i>	
1634	<i>ference on computer vision and pattern recognition</i> ,	
1635	pages 8317–8326.	
1636	Shezheng Song, Xiaopeng Li, and Shasha Li. 2023.	
1637	How to Bridge the Gap between Modalities: A Com-	
1638	prehensive Survey on Multimodal Large Language	
1639	Model. <i>arXiv preprint arXiv:2311.07594</i> .	
	Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma,	1640
	Abhishek Kumar, Stefano Ermon, and Ben Poole.	1641
	2021. Score-Based Generative Modeling through	1642
	Stochastic Differential Equations. In <i>International</i>	1643
	<i>Conference on Learning Representations</i> .	1644
	Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan	1645
	Wang, and Deng Cai. 2023. Pandagpt: One	1646
	model to instruction-follow them all. <i>arXiv preprint</i>	1647
	<i>arXiv:2305.16355</i> .	1648
	Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang,	1649
	Qiyong Yu, Zhengxiong Luo, Yueze Wang, Yongming	1650
	Rao, Jingjing Liu, Tiejun Huang, et al. 2023a. Gen-	1651
	erative multimodal models are in-context learners.	1652
	<i>arXiv preprint arXiv:2312.13286</i> .	1653
	Quan Sun, Qiyong Yu, Yufeng Cui, Fan Zhang,	1654
	Xiaosong Zhang, Yueze Wang, Hongcheng Gao,	1655
	Jingjing Liu, Tiejun Huang, and Xinlong Wang. 2024.	1656
	Generative pretraining in multimodality. In <i>The</i>	1657
	<i>Twelfth International Conference on Learning Repre-</i>	1658
	<i>sentations</i> .	1659
	Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu,	1660
	Chunyu Li, Yikang Shen, Chuang Gan, Liang-Yan	1661
	Gui, Yu-Xiong Wang, Yiming Yang, et al. 2023b.	1662
	Aligning large multimodal models with factually aug-	1663
	mented rlhf. <i>arXiv preprint arXiv:2309.14525</i> .	1664
	Dídac Surís, Sachit Menon, and Carl Vondrick. 2023.	1665
	ViperGPT: Visual inference via python execution for	1666
	reasoning. <i>arXiv preprint arXiv:2303.08128</i> .	1667
	Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao	1668
	Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and Chao	1669
	Zhang. 2023a. Salmonn: Towards generic hearing	1670
	abilities for large language models. <i>arXiv preprint</i>	1671
	<i>arXiv:2310.13289</i> .	1672
	Zineng Tang, Ziyi Yang, Mahmoud Khademi, Yang Liu,	1673
	Chenguang Zhu, and Mohit Bansal. 2023b. CoDi-2:	1674
	In-Context, Interleaved, and Interactive Any-to-Any	1675
	Generation. <i>arXiv preprint arXiv:2311.18775</i> .	1676
	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng,	1677
	and Mohit Bansal. 2023c. Any-to-Any Generation	1678
	via Composable Diffusion. In <i>Thirty-seventh Confer-</i>	1679
	<i>ence on Neural Information Processing Systems</i> .	1680
	Yi Tay, Mostafa Dehghani, Vinh Q Tran, Xavier Gar-	1681
	cia, Jason Wei, Xuezhi Wang, Hyung Won Chung,	1682
	Dara Bahri, Tal Schuster, Steven Zheng, et al. 2022.	1683
	UI2: Unifying language learning paradigms. In <i>The</i>	1684
	<i>Eleventh International Conference on Learning Repre-</i>	1685
	<i>sentations</i> .	1686
	Gemini Team, Rohan Anil, Sebastian Borgeaud,	1687
	Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,	1688
	Radu Soricut, Johan Schalkwyk, Andrew M Dai,	1689
	Anja Hauth, et al. 2023. Gemini: a family of	1690
	highly capable multimodal models. <i>arXiv preprint</i>	1691
	<i>arXiv:2312.11805</i> .	1692
	InternLM Team. 2023. Internlm: A multilingual lan-	1693
	guage model with progressively enhanced capabili-	1694
	ties.	1695

1696	Yi Team. 2023. <a href="#">Yi-VL</a> .	understanding of the open world. <i>arXiv preprint arXiv:2308.01907</i> .	1752 1753
1697	Changyao Tian, Xizhou Zhu, Yuwen Xiong, Weiyun Wang, Zhe Chen, Wenhai Wang, Yuntao Chen, Lewei Lu, Tong Lu, Jie Zhou, et al. 2024. MM-Interleaved: Interleaved Image-Text Generative Modeling via Multi-modal Feature Synchronizer. <i>arXiv preprint arXiv:2401.10208</i> .	Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. 2022b. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. <i>arXiv preprint arXiv:2208.10442</i> .	1754 1755 1756 1757 1758 1759
1703	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. 2023d. Image as a Foreign Language: BEiT Pretraining for Vision and Vision-Language Tasks. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 19175–19186.	1760 1761 1762 1763 1764 1765 1766 1767
1709	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	Xinyu Wang, Bohan Zhuang, and Qi Wu. 2024c. ModaVerse: Efficiently Transforming Modalities with LLMs. <i>arXiv preprint arXiv:2401.06395</i> .	1768 1769 1770
1715	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, En Yu, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. 2024. Small Language Model Meets with Reinforced Vision Vocabulary. <i>arXiv preprint arXiv:2401.12503</i> .	1771 1772 1773 1774 1775
1720	Chenyu Wang, Weixin Luo, Qianyu Chen, Haonan Mai, Jindi Guo, Sixun Dong, Zhengxin Li, Lin Ma, Shenghua Gao, et al. 2024a. Tool-LMM: A Large Multi-Modal Model for Tool Agent Learning. <i>arXiv preprint arXiv:2401.10727</i> .	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned Language Models are Zero-Shot Learners. In <i>International Conference on Learning Representations</i> .	1776 1777 1778 1779 1780
1725	Dongsheng Wang, Natraj Raman, Mathieu Sibue, Zhiqiang Ma, Petr Babkin, Simerjot Kaur, Yulong Pei, Armineh Nourbakhsh, and Xiaomo Liu. 2023a. DocLLM: A layout-aware generative language model for multimodal document understanding. <i>arXiv preprint arXiv:2401.00908</i> .	Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. 2023a. Visual chatgpt: Talking, drawing and editing with visual foundation models. <i>arXiv preprint arXiv:2303.04671</i> .	1781 1782 1783 1784 1785
1731	Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. 2024b. Mobile-Agent: Autonomous Multi-Modal Mobile Device Agent with Visual Perception. <i>arXiv preprint arXiv:2401.16158</i> .	Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Chunyi Li, Wenxiu Sun, Qiong Yan, Guangtao Zhai, et al. 2023b. Q-bench: A benchmark for general-purpose foundation models on low-level vision. <i>arXiv preprint arXiv:2309.14181</i> .	1786 1787 1788 1789 1790 1791
1736	Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022a. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In <i>International Conference on Machine Learning</i> , pages 23318–23340. PMLR.	Jiahong Wu, He Zheng, Bo Zhao, Yixin Li, Baoming Yan, Rui Liang, Wenjia Wang, Shipei Zhou, Guosen Lin, Yanwei Fu, et al. 2017. Ai challenger: A large-scale dataset for going deeper in image understanding. <i>arXiv preprint arXiv:1711.06475</i> .	1792 1793 1794 1795 1796
1743	Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. 2023b. CogVLM: Visual expert for pretrained language models. <i>arXiv preprint arXiv:2311.03079</i> .	Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and Philip S Yu. 2023c. Multimodal large language models: A survey. <i>arXiv preprint arXiv:2311.13165</i> .	1797 1798 1799 1800
1748	Weiyun Wang, Min Shi, Qingyun Li, Wenhai Wang, Zhenhang Huang, Linjie Xing, Zhe Chen, Hao Li, Xizhou Zhu, Zhiguo Cao, et al. 2023c. The all-seeing project: Towards panoptic visual recognition and	Penghao Wu and Saining Xie. 2023. V*: Guided Visual Search as a Core Mechanism in Multimodal LLMs. <i>arXiv preprint arXiv:2312.14135</i> , 17.	1801 1802 1803
1749		Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2023d. Next-gpt: Any-to-any multimodal llm. <i>arXiv preprint arXiv:2309.05519</i> .	1804 1805 1806

1807	Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan,	Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Haowei	1864
1808	Thuy-Trang Vu, and Gholamreza Haffari. 2024. Con-	Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou.	1865
1809	tinual Learning for Large Language Models: A Sur-	2023c. mplug-owl2: Revolutionizing multi-modal	1866
1810	vey. <i>arXiv preprint arXiv:2402.01364</i> .	large language model with modality collaboration.	1867
1811	Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Tay-	<i>arXiv preprint arXiv:2311.04257</i> .	1868
1812	lor Berg-Kirkpatrick, and Shlomo Dubnov. 2023e.	Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun,	1869
1813	Large-scale contrastive language-audio pretraining	Tong Xu, and Enhong Chen. 2023a. A Survey on	1870
1814	with feature fusion and keyword-to-caption augmen-	Multimodal Large Language Models. <i>arXiv preprint</i>	1871
1815	tation. In <i>ICASSP 2023-2023 IEEE International</i>	<i>arXiv:2306.13549</i> .	1872
1816	<i>Conference on Acoustics, Speech and Signal Process-</i>	Zhenfei Yin, Jiong Wang, Jianjian Cao, Zhelun Shi,	1873
1817	<i>ing (ICASSP)</i> , pages 1–5. IEEE.	Dingning Liu, Mukai Li, Lu Sheng, Lei Bai, Xi-	1874
1818	Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. Msr-	aoshui Huang, Zhiyong Wang, et al. 2023b. Lamm:	1875
1819	vtv: A large video description dataset for bridging	Language-assisted multi-modal instruction-tuning	1876
1820	video and language. In <i>Proceedings of the IEEE con-</i>	dataset, framework, and benchmark. <i>arXiv preprint</i>	1877
1821	<i>ference on computer vision and pattern recognition</i> ,	<i>arXiv:2306.06687</i> .	1878
1822	pages 5288–5296.	Peter Young, Alice Lai, Micah Hodosh, and Julia Hock-	1879
1823	Rui Yan, Mike Zheng Shou, Yixiao Ge, Alex Jinpeng	enmaier. 2014. From image descriptions to visual	1880
1824	Wang, Xudong Lin, Guanyu Cai, and Jinhui Tang.	denotations: New similarity metrics for semantic	1881
1825	2021. Video-text pre-training with learned regions.	inference over event descriptions. <i>Transactions of the</i>	1882
1826	<i>arXiv preprint arXiv:2112.01194</i> .	<i>Association for Computational Linguistics</i> , 2:67–78.	1883
1827	Siming Yan, Min Bai, Weifeng Chen, Xiong Zhou, Qix-	Licheng Yu, Patrick Poirson, Shan Yang, Alexander C	1884
1828	ing Huang, and Li Erran Li. 2024. ViGoR: Improving	Berg, and Tamara L Berg. 2016. Modeling context	1885
1829	Visual Grounding of Large Vision Language Models	in referring expressions. In <i>Computer Vision–ECCV</i>	1886
1830	with Fine-Grained Reward Modeling. <i>arXiv preprint</i>	<i>2016: 14th European Conference, Amsterdam, The</i>	1887
1831	<i>arXiv:2402.06118</i> .	<i>Netherlands, October 11-14, 2016, Proceedings, Part</i>	1888
1832	Jinyu Yang, Jiali Duan, Son Tran, Yi Xu, Sampath	<i>II 14</i> , pages 69–85. Springer.	1889
1833	Chanda, Liqun Chen, Belinda Zeng, Trishul Chilimbi,	Lili Yu, Bowen Shi, Ramakanth Pasunuru, Benjamin	1890
1834	and Junzhou Huang. 2022. Vision-language pre-	Muller, Olga Golovneva, Tianlu Wang, Arun Babu,	1891
1835	training with triple contrastive learning. In <i>Proceed-</i>	Binh Tang, Brian Karrer, Shelly Sheynin, et al.	1892
1836	<i>ings of the IEEE/CVF Conference on Computer Vi-</i>	2023a. Scaling autoregressive multi-modal models:	1893
1837	<i>sion and Pattern Recognition</i> , pages 15671–15680.	Pretraining and instruction tuning. <i>arXiv preprint</i>	1894
1838	Ling Yang, Zhaochen Yu, Chenlin Meng, Minkai Xu,	<i>arXiv:2309.02591</i> .	1895
1839	Stefano Ermon, and Bin Cui. 2024. Mastering Text-	Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwan He, Yifeng	1896
1840	to-Image Diffusion: Recaptioning, Planning, and	Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao	1897
1841	Generating with Multimodal LLMs. <i>arXiv preprint</i>	Zheng, Maosong Sun, et al. 2023b. RLhf-v: Towards	1898
1842	<i>arXiv:2401.11708</i> .	trustworthy mlms via behavior alignment from fine-	1899
1843	Zhen Yang, Yingxue Zhang, Fandong Meng, and Jie	grained correctional human feedback. <i>arXiv preprint</i>	1900
1844	Zhou. 2023a. TEAL: Tokenize and Embed ALL for	<i>arXiv:2312.00849</i> .	1901
1845	Multi-modal Large Language Models. <i>arXiv preprint</i>	Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang,	1902
1846	<i>arXiv:2311.04589</i> .	Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan	1903
1847	Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin	Wang. 2023c. Mm-vet: Evaluating large multimodal	1904
1848	Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu,	models for integrated capabilities. <i>arXiv preprint</i>	1905
1849	Ce Liu, Michael Zeng, and Lijuan Wang. 2023b.	<i>arXiv:2308.02490</i> .	1906
1850	Mm-react: Prompting chatgpt for multimodal rea-	Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang,	1907
1851	soning and action. <i>arXiv preprint arXiv:2303.11381</i> .	Jie Zhou, and Jiwen Lu. 2022. Point-bert: Pre-	1908
1852	Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye,	training 3d point cloud transformers with masked	1909
1853	Ming Yan, Yuhao Dan, Chenlin Zhao, Guohai Xu,	point modeling. In <i>Proceedings of the IEEE/CVF</i>	1910
1854	Chenliang Li, Junfeng Tian, et al. 2023a. mplug-	<i>Conference on Computer Vision and Pattern Recog-</i>	1911
1855	docowl: Modularized multimodal large language	<i>nition</i> , pages 19313–19322.	1912
1856	model for document understanding. <i>arXiv preprint</i>	Yuqian Yuan, Wentong Li, Jian Liu, Dongqi Tang, Xin-	1913
1857	<i>arXiv:2307.02499</i> .	jie Luo, Chi Qin, Lei Zhang, and Jianke Zhu. 2023a.	1914
1858	Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye,	Osprey: Pixel Understanding with Visual Instruction	1915
1859	Ming Yan, Yiyang Zhou, Junyang Wang, An-	Tuning. <i>arXiv preprint arXiv:2312.10032</i> .	1916
1860	wen Hu, Pengcheng Shi, Yaya Shi, et al. 2023b.	Zhengqing Yuan, Zhaoxu Li, and Lichao Sun. 2023b.	1917
1861	mplug-owl: Modularization empowers large lan-	TinyGPT-V: Efficient Multimodal Large Language	1918
1862	guage models with multimodality. <i>arXiv preprint</i>	Model via Small Backbones. <i>arXiv preprint</i>	1919
1863	<i>arXiv:2304.14178</i> .	<i>arXiv:2312.16862</i> .	1920



1921	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2023. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. <i>arXiv preprint arXiv:2311.16502</i> .	Frontiers in Spiking Neural Networks. In <i>Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022</i> , pages 5670–5677.	1976 1977 1978 1979
1927	Rowan Zellers, Jiasen Lu, Ximing Lu, Youngjae Yu, Yanpeng Zhao, Mohammadreza Salehi, Aditya Kusupati, Jack Hessel, Ali Farhadi, and Yejin Choi. 2022. Merlot reserve: Neural script knowledge through vision and language and sound. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 16375–16387.	Ge Zhang, Xinrun Du, Bei Chen, Yiming Liang, Tongxu Luo, Tianyu Zheng, Kang Zhu, Yuyang Cheng, Chunpu Xu, Shuyue Guo, et al. 2024. CM-MMU: A Chinese Massive Multi-discipline Multimodal Understanding Benchmark. <i>arXiv preprint arXiv:2401.11944</i> .	1980 1981 1982 1983 1984 1985
1934	Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022a. GLM-130B: An Open Bilingual Pre-trained Model. In <i>The Eleventh International Conference on Learning Representations</i> .	Hang Zhang, Xin Li, and Lidong Bing. 2023e. Video-LLaMA: An Instruction-tuned Audio-Visual Language Model for Video Understanding. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023 - System Demonstrations, Singapore, December 6-10, 2023</i> , pages 543–553.	1986 1987 1988 1989 1990 1991 1992
1940	Yan Zeng, Hanbo Zhang, Jiani Zheng, Jiangnan Xia, Guoqiang Wei, Yang Wei, Yuchen Zhang, and Tao Kong. 2023. What Matters in Training a GPT4-Style Language Model with Multimodal Inputs? <i>arXiv preprint arXiv:2307.02469</i> .	Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel Ni, and Heung-Yeung Shum. 2022b. DINO: DETR with Improved DeNoising Anchor Boxes for End-to-End Object Detection. In <i>The Eleventh International Conference on Learning Representations</i> .	1993 1994 1995 1996 1997 1998
1945	Yan Zeng, Xinsong Zhang, and Hang Li. 2022b. Multi-Grained Vision Language Pre-Training: Aligning Texts with Visual Concepts. In <i>International Conference on Machine Learning</i> , pages 25994–26009. PMLR.	Jeffrey O Zhang, Alexander Sax, Amir Zamir, Leonidas Guibas, and Jitendra Malik. 2020. Side-tuning: a baseline for network adaptation via additive side networks. In <i>Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16</i> , pages 698–714. Springer.	1999 2000 2001 2002 2003 2004
1950	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. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 15757–15773.	Pan Zhang, Xiaoyi Dong Bin Wang, Yuhang Cao, Chao Xu, Linke Ouyang, Zhiyuan Zhao, Shuangrui Ding, Songyang Zhang, Haodong Duan, Hang Yan, et al. 2023f. Internlm-xcomposer: A vision-language large model for advanced text-image comprehension and composition. <i>arXiv preprint arXiv:2309.15112</i> .	2005 2006 2007 2008 2009 2010
1957	Duzhen Zhang, Wei Cong, Jiahua Dong, Yahan Yu, Xiuyi Chen, Yonggang Zhang, and Zhen Fang. 2023b. Continual Named Entity Recognition without Catastrophic Forgetting. In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Kai Chen, and Ping Luo. 2023g. Gpt4roi: Instruction tuning large language model on region-of-interest. <i>arXiv preprint arXiv:2307.03601</i> .	2011 2012 2013 2014 2015
1962	Duzhen Zhang, Hongliu Li, Wei Cong, Rongtao Xu, Jiahua Dong, and Xiuyi Chen. 2023c. Task relation distillation and prototypical pseudo label for incremental named entity recognition. In <i>Proceedings of the 32nd ACM International Conference on Information and Knowledge Management</i> , pages 3319–3329.	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022c. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> .	2016 2017 2018 2019 2020
1968	Duzhen Zhang, Yahan Yu, Feilong Chen, and Xiuyi Chen. 2023d. Decomposing Logits Distillation for Incremental Named Entity Recognition. In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 1919–1923.	Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. 2023h. Llavav: Enhanced visual instruction tuning for text-rich image understanding. <i>arXiv preprint arXiv:2306.17107</i> .	2021 2022 2023 2024 2025
1974	Duzhen Zhang, Tielin Zhang, Shuncheng Jia, Qingyu Wang, and Bo Xu. 2022a. Recent Advances and New	Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li, Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin, Tong Luo, Yaqian Li, Shilong Liu, et al. 2023i. Recognize Anything: A Strong Image Tagging Model. <i>arXiv preprint arXiv:2306.03514</i> .	2026 2027 2028 2029 2030

2031	Bingchen Zhao, Haoqin Tu, Chen Wei, and Cihang Xie.	Dawei Yin. 2024b. VisLingInstruct: Elevating Zero-Shot Learning in Multi-Modal Language Models with Autonomous Instruction Optimization. <i>arXiv preprint arXiv:2402.07398</i> .	2085
2032	2024. Tuning LayerNorm in Attention: Towards Efficient Multimodal LLM Finetuning. In <i>The Twelfth International Conference on Learning Representations</i> .		2086
2033			2087
2034			2088
2035			
2036	Bo Zhao, Boya Wu, and Tiejun Huang. 2023a. Svit: Scaling up visual instruction tuning. <i>arXiv preprint arXiv:2307.04087</i> .	Jinguo Zhu, Xiaohan Ding, Yixiao Ge, Yuying Ge, Sijie Zhao, Hengshuang Zhao, Xiaohua Wang, and Ying Shan. 2023b. V1-gpt: A generative pre-trained transformer for vision and language understanding and generation. <i>arXiv preprint arXiv:2312.09251</i> .	2089
2037			2090
2038			2091
2039	Liang Zhao, En Yu, Zheng Ge, Jinrong Yang, Hao-ran Wei, Hongyu Zhou, Jianjian Sun, Yuang Peng, Runpei Dong, Chunrui Han, et al. 2023b. Chatspot: Bootstrapping multimodal llms via precise referring instruction tuning. <i>arXiv preprint arXiv:2307.09474</i> .		2092
2040			2093
2041		Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. 2023c. Multimodal c4: An open, billion-scale corpus of images interleaved with text. <i>arXiv preprint arXiv:2304.06939</i> .	2094
2042			2095
2043			2096
2044	Min Zhao, Fan Bao, Chongxuan Li, and Jun Zhu. 2022. EGSD: Unpaired Image-to-Image Translation via Energy-Guided Stochastic Differential Equations. In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .		2097
2045			2098
2046			2099
2047		Yichen Zhu, Minjie Zhu, Ning Liu, Zhicai Ou, Xiaofeng Mou, and Jian Tang. 2024c. LLaVA-Phi: Efficient Multi-Modal Assistant with Small Language Model. <i>arXiv preprint arXiv:2401.02330</i> .	2100
2048			2101
2049			2102
2050			2103
2051	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023c. A survey of large language models. <i>arXiv preprint arXiv:2303.18223</i> .	Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. 2016. Visual7w: Grounded question answering in images. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 4995–5004.	2104
2052			2105
2053			2106
2054			2107
2055			2108
2056	Yang Zhao, Zhijie Lin, Daquan Zhou, Zilong Huang, Jiashi Feng, and Bingyi Kang. 2023d. Bubogpt: Enabling visual grounding in multi-modal llms. <i>arXiv preprint arXiv:2307.08581</i> .	Bojia Zi, Xianbiao Qi, Lingzhi Wang, Jianan Wang, Kam-Fai Wong, and Lei Zhang. 2023. Delta-lora: Fine-tuning high-rank parameters with the delta of low-rank matrices. <i>arXiv preprint arXiv:2309.02411</i> .	2109
2057			2110
2058			2111
2059			2112
2060	Junhao Zheng, Qianli Ma, Zhen Liu, Binquan Wu, and Huawen Feng. 2024. Beyond Anti-Forgetting: Multimodal Continual Instruction Tuning with Positive Forward Transfer. <i>arXiv preprint arXiv:2401.09181</i> .	Yongshuo Zong, Ondrej Bohdal, Tingyang Yu, Yongxin Yang, and Timothy Hospedales. 2024. Safety Fine-Tuning at (Almost) No Cost: A Baseline for Vision Large Language Models. <i>arXiv preprint arXiv:2402.02207</i> .	2113
2061			2114
2062			2115
2063			2116
2064	Junhao Zheng, Shengjie Qiu, and Qianli Ma. 2023a. Learn or Recall? Revisiting Incremental Learning with Pre-trained Language Models. <i>arXiv preprint arXiv:2312.07887</i> .		2117
2065			
2066			
2067			
2068	Kaizhi Zheng, Xuehai He, and Xin Eric Wang. 2023b. Minigpt-5: Interleaved vision-and-language generation via generative vokens. <i>arXiv preprint arXiv:2310.02239</i> .		
2069			
2070			
2071			
2072	Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiayi Cui, HongFa Wang, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, et al. 2024a. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment. In <i>The Twelfth International Conference on Learning Representations</i> .		
2073			
2074			
2075			
2076			
2077			
2078			
2079	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023a. Minigpt-4: Enhancing vision-language understanding with advanced large language models. <i>arXiv preprint arXiv:2304.10592</i> .		
2080			
2081			
2082			
2083	Dongsheng Zhu, Xunzhu Tang, Weidong Han, Jinghui Lu, Yukun Zhao, Guoliang Xing, Junfeng Wang, and		
2084			

2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
  
2148  
2149  
2150  
2151  
2152  
2153  
  
2154  
2155  
2156  
2157  
  
2158  
2159  
2160  
2161  
2162  
2163

## A Related Surveys

Prior to the emergence of LLMs, several surveys on traditional MM PT have been conducted (Ruan and Jin, 2022; Du et al., 2022a; Long et al., 2022; Chen et al., 2023a). Most of these models entail a substantial computational cost during the PT phase, attributable to end-to-end training using large-scale models and datasets. As a consequence of not incorporating LLMs, these models suffer from deficiencies in instruction following, ICL, CoT, and interactive capabilities. Moreover, the training pipeline solely encompasses the PT phase without the inclusion of an IT stage.

In recent times, several surveys have emerged on MM-LLMs. Yin et al. and Wu et al. exclusively delve into early VL understanding models. Huang et al. place a primary emphasis on visual IT, while Song et al. focus on modal alignment methods. Lastly, Cui et al. provide a comprehensive review of the applications of MM-LLMs within the realm of autonomous driving.

Compared with their works, the main distinctions are outlined as follows:

- We have comprehensively covered nearly all MM-LLMs over the past year, totaling around 120 or more, including not only understanding models but also generative models. Our coverage extends beyond VL modalities to encompass various modes such as audio and 3D point cloud;
- To offer readers a comprehensive understanding of MM-LLMs, we have introduced a general model architecture that incorporates any-to-any modality transformations, offering a detailed overview of the functional roles and implementation choices for each component;
- We have summarized the developmental trends of existing MM-LLMs and provided some training recipes that can enhance effectiveness;
- We have established an open-source website for MM-LLMs researchers, supporting crowd-sourced updates and aiming to facilitate collaboration in the MM-LLMs field. We anticipate that this survey will illuminate future research in the MM-LLMs domain.

## B Modality Encoder

In the following, we provide a brief introduction to some mainstream modality encoders.

### B.1 Visual Modality

**NFNet-F6** (Brock et al., 2021) is a normalizer-free ResNet (He et al., 2016), showcasing an adaptive gradient clipping that allows training on extensively augmented datasets while achieving SOTA levels of image recognition.

**ViT** (Dosovitskiy et al., 2020) applies the Transformer (Vaswani et al., 2017) to images by first dividing the image into patches. It then undergoes linear projection to flatten the patches, followed by encoding via Transformer blocks.

**CLIP ViT** (Radford et al., 2021) builds connections between text and images, comprising a ViT and a text encoder. With a vast amount of text-image pairs, it optimizes ViT by contrastive learning, treating paired text and images as positive samples and others as negative ones.

**Eva-CLIP ViT** (Fang et al., 2023) stabilizes the training and optimization process of the massive CLIP, offering new directions in expanding and accelerating the expensive training of MM base models.

### B.2 Audio Modality

**C-Former** (Chen et al., 2023b) employs the CIF (Dong and Xu, 2020; Zhang et al., 2022a; Han et al., 2022, 2023) for sequence transduction and a Transformer to extract audio features.

**HuBERT** (Hsu et al., 2021) is a self-supervised speech representation learning framework based on BERT (Kenton and Toutanova, 2019), achieved by the masked prediction of discrete hidden units. It has the capability to convert continuous speech signals into a sequence of discrete units.

**BEATs** (Chen et al., 2023g) is an iterative audio pre-training framework designed to learn Bidirectional Encoder representations from Audio Transformers.

## C Mainstream PEFT Methods

PEFT entails maintaining the pre-trained LLM in a frozen state while adjusting a small number of additional trainable parameters. In the following section, we revisit several representative PEFT methods, where  $x$  and  $h$  represent the input and output

2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209

of the original module, and  $h'$  signifies the output of this module when attached with PEFT.

**Prefix-tuning** (Li and Liang, 2021; Lester et al., 2021) involves the addition of learnable tokens to the keys and values of the attention module. This process is formulated as follows:

$$h' = \text{Attn}(\mathbf{x}\mathbf{W}_q, [\mathbf{P}_k, \mathbf{x}\mathbf{W}_k], [\mathbf{P}_v, \mathbf{x}\mathbf{W}_v]), \quad (6)$$

with  $\mathbf{P}_k, \mathbf{P}_v \in \mathbb{R}^{l \times d}$  representing two sets of prefix tokens.  $[\cdot, \cdot]$  denotes concatenation, and  $\text{Attn}$  is defined as:

$$\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) := \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}.$$

**Adapter** (Houlsby et al., 2019; He et al., 2021; Rebuffi et al., 2017; Zhang et al., 2020) is typically a residual block consisting of a down-projection matrix  $\mathbf{A}$ , a nonlinear activation function  $\sigma(\cdot)$ , and an up-projection matrix  $\mathbf{B}$ . It can be inserted into any layer of the pre-trained LLM, formulated as follows:

$$h' = h + \sigma(\mathbf{x}\mathbf{A})\mathbf{B}. \quad (7)$$

**LoRA** (Hu et al., 2021) is the most commonly used PEFT method. It assumes that the change in parameters occurs within a low-rank space. Given a pre-trained matrix  $\mathbf{W} \in \mathbb{R}^{c \times d}$ , LoRA learns an incremental update  $\Delta\mathbf{W}$  and decomposes  $\Delta\mathbf{W}$  into a matrix multiplication between two low-rank matrices  $\mathbf{A} \in \mathbb{R}^{c \times r}$  and  $\mathbf{B} \in \mathbb{R}^{r \times d}$ , where  $r \ll \min(c, d)$ . LoRA follows the forward process as outlined below:

$$h = \mathbf{W}\mathbf{x} + \Delta\mathbf{W}\mathbf{x} = \mathbf{W}\mathbf{x} + \mathbf{A}\mathbf{B}\mathbf{x}. \quad (8)$$

**QLoRA** (Dettmers et al., 2023) is a quantized LoRA. The underlying principle of QLoRA includes the quantization of pre-trained weights to 4 bits, followed by the execution of PEFT using LoRA.

**LayerNorm tuning** (Zhao et al., 2024) presents an efficient strategy to transform LLMs into MM-LLMs, which tunes LayerNorm in attention block yielding strong MM performance compared with full parameter finetuning or LoRA.

In addition to the aforementioned PEFT methods, there are several others, including P-tuning (Liu et al., 2022), P-tuning v2 (Liu et al., 2021a), Adapt-Bias (Fu et al., 2022), Compacter (Karimi Mahabadi et al., 2021), AdapterFormer (Chen et al., 2022a), XTuner (Contributors, 2023), P-LoRA (Dong et al., 2024b), MoLE (Chen et al., 2024), and Delta-LoRA (Zi et al., 2023).

## D Representative LLMs

The representative LLM Backbones in existing MM-LLMs research are as follows:

**Flan-T5** (Chung et al., 2022) investigates IT for T5 (Raffel et al., 2020), an encoder-decoder architecture using unified text-to-text training for all natural language processing issues, exhibiting robust zero-shot and CoT capabilities.

**ChatGLM** is a Chinese-English bilingual dialogue model,<sup>2</sup> optimized by an auto-regressive mask infilling objective. It is based on the GLM (Du et al., 2022b; Zeng et al., 2022a) architecture, optimized for Chinese question answering and dialogues.

**InternLM** (Team, 2023) is a multilingual trillion-parameter foundation model trained on over a trillion tokens of data. Based on this foundation, the model utilizes high-quality human-annotated dialogue data combined with RLHF to respond to complex instructions during human interactions, exhibiting responses that align with human ethics and values.

**UL2** (Tay et al., 2022) is an encoder-decoder model trained utilizing a mixture of denoisers objectives, surpassing T5 on numerous benchmarks.

**Qwen** (Bai et al., 2023a) is trained on large-scale and diverse datasets, with a primary focus on Chinese and English. It employs SFT and RLHF techniques for alignment, resulting in dialogue models like Qwen-Chat.

**Chinchilla** (Hoffmann et al., 2022) is a causal decoder, trained on extensive text data. It posits that model size should double for every doubling of training tokens.

**OPT** (Zhang et al., 2022c) is a GPT-3 (Brown et al., 2020) clone, striving to release an open-source model that replicates the performance of GPT-3.

**PaLM** (Chowdhery et al., 2023) is a causal decoder structure with parallel attention and feed-forward layers, enabling training speeds up to 15 times faster. Notable changes contain RoPE embeddings, SwiGLU activation, multi-query attention, and etc.

<sup>2</sup><https://github.com/THUDM/ChatGLM-6B>



2297	<b>LLaMA</b> (Touvron et al., 2023a) comprises	<b>VideoChat</b> (Li et al., 2023f) pioneers an efficient	2344
2298	decoder-only models with efficient causal atten-	chat-centric MM-LLM for video understanding di-	2345
2299	tion.	alogue, setting standards for future research in this	2346
2300	<b>LLaMA-2</b> (Touvron et al., 2023b) focuses on	domain and offering protocols for both academia	2347
2301	fine-tuning a superior and safer LLaMA-2-Chat	and industry.	2348
2302	model for conversation generation, incorporating	<b>InstructBLIP</b> (Dai et al., 2023) is trained based	2349
2303	40% more training data with grouped-query atten-	on the pre-trained BLIP-2 model, updating only	2350
2304	tion and a larger context length.	the Q-Former during MM IT. By introducing	2351
2305	<b>Vicuna</b> (Chiang et al., 2023) is a model built	instruction-aware visual feature extraction and cor-	2352
2306	on top of LLaMA, utilizing user dialogue data ob-	responding instructions, the model enables the ex-	2353
2307	tained from ShareGPT.com and trained by SFT.	traction of flexible and diverse features.	2354
2308	<b>E SOTA MM-LLMs</b>	<b>PandaGPT</b> (Su et al., 2023) is a pioneering	2355
2309	In the following, we will provide a brief introduc-	general-purpose model with the capability to com-	2356
2310	tion to the core contributions of some representa-	prehend and act upon instructions across 6 differ-	2357
2311	MM-LLMs.	ent modalities: text, image/video, audio, thermal,	2358
2312	<b>Flamingo</b> (Alayrac et al., 2022) represents a se-	depth, and inertial measurement units.	2359
2313	ries of Visual Language (VL) Models designed for	<b>(PaLI-X</b> (Chen et al., 2023h) is trained using	2360
2314	processing interleaved visual data and text, gener-	mixed VL objectives and unimodal objectives, in-	2361
2315	ating free-form text as the output.	cluding prefix completion and masked-token com-	2362
2316	<b>BLIP-2</b> (Li et al., 2023e) introduces a more	pletion. This approach proves effective for both	2363
2317	resource-efficient framework, comprising the	downstream task results and achieving the Pareto	2364
2318	lightweight Q-Former to bridge modality gaps and	frontier in the fine-tuning setting.	2365
2319	the utilization of frozen LLMs. Leveraging LLMs,	<b>Video-LLaMA</b> (Zhang et al., 2023e) introduces	2366
2320	BLIP-2 can be guided for zero-shot image-to-text	a multi-branch cross-modal PT framework, en-	2367
2321	generation using natural language prompts.	abling LLMs to simultaneously process the vision	2368
2322	<b>LLaVA</b> (Liu et al., 2023e) pioneers the trans-	and audio content of a given video while engag-	2369
2323	fer of IT techniques to the MM domain. Ad-	ing in conversations with humans. This framework	2370
2324	dressing data scarcity, LLaVA introduces a novel	aligns vision with language as well as audio with	2371
2325	open-source MM instruction-following dataset cre-	language.	2372
2326	ated using ChatGPT/GPT-4, alongside the MM	<b>Video-ChatGPT</b> (Maaz et al., 2023) is a model	2373
2327	instruction-following benchmark, LLaVA-Bench.	specifically designed for video conversations, ca-	2374
2328	<b>MiniGPT-4</b> (Zhu et al., 2023a) proposes a	pable of generating discussions about videos by	2375
2329	streamlined approach where training only one lin-	integrating spatiotemporal vision representations.	2376
2330	ear layer aligns the pre-trained vision encoder with	<b>Shikra</b> (Chen et al., 2023e) introduces a sim-	2377
2331	the LLM. This efficient method enables the repli-	ple and unified pre-trained MM-LLM tailored for	2378
2332	cation of the exhibited capabilities of GPT-4.	Referential Dialogue, a task involving discussions	2379
2333	<b>mPLUG-Owl</b> (Ye et al., 2023b) presents a novel	about regions and objects in images. This model	2380
2334	modularized training framework for MM-LLMs,	demonstrates commendable generalization ability,	2381
2335	incorporating the visual context. To assess differ-	effectively handling unseen settings.	2382
2336	ent models' performance in MM tasks, the frame-	<b>DLP</b> (Jian et al., 2023) proposes the P-Former	2383
2337	work includes an instructional evaluation dataset called	to predict the ideal prompt, trained on a dataset	2384
2338	OwIEval.	of single-modal sentences. This showcases the	2385
2339	<b>X-LLM</b> (Chen et al., 2023b) is expanded to var-	feasibility of single-modal training to enhance MM	2386
2340	ious modalities, including audio, and demonstrates	learning.	2387
2341	strong scalability. Leveraging the language trans-	<b>BuboGPT</b> (Zhao et al., 2023d) is a model con-	2388
2342	ferability of the Q-Former, X-LLM is successfully	structed by learning a shared semantic space for a	2389
2343	applied in the context of Sino-Tibetan Chinese.	comprehensive understanding of MM content. It	2390

2391	explores fine-grained relationships among different modalities such as image, text, and audio.	
2392		
2393	<b>ChatSpot</b> (Zhao et al., 2023b) introduces a simple yet potent method for finely adjusting precise referring instructions for MM-LLM, facilitating fine-grained interactions. The incorporation of precise referring instructions, consisting of image- and region-level instructions, enhances the integration of multi-grained VL task descriptions.	2439
2394		2440
2395		2441
2396		2442
2397		2443
2398		2444
2399		2445
2400	<b>Qwen-VL</b> (Bai et al., 2023b) is a multi-lingual MM-LLM that supports both English and Chinese. Qwen-VL also allows the input of multiple images during the training phase, improving its ability to understand the vision context.	
2401		
2402		
2403		
2404		
2405	<b>NExT-GPT</b> (Wu et al., 2023d) is an end-to-end, general-purpose any-to-any MM-LLM that supports the free input and output of image, video, audio, and text. It employs a lightweight alignment strategy, utilizing LLM-centric alignment in the encoding phase and instruction-following alignment in the decoding phase.	2446
2406		2447
2407		2448
2408		2449
2409		2450
2410		2451
2411		2452
2412	<b>MiniGPT-5</b> (Zheng et al., 2023b) is an MM-LLM integrated with inversion to generative tokens and integration with Stable Diffusion. It excels in performing interleaved VL outputs for MM generation. The inclusion of classifier-free guidance during the training phase enhances the quality of generation.	2453
2413		2454
2414		2455
2415		2456
2416		2457
2417		2458
2418		2459
2419	<b>LLaVA-1.5</b> (Liu et al., 2023d) reports simple modifications to the LLaVA framework, including applying an MLP projection and introducing VQA data tailored for academic tasks, along with simple response formatting prompts. These adjustments result in enhanced capabilities for MM understanding.	2460
2420		2461
2421		2462
2422		2463
2423		2464
2424		2465
2425		2466
2426	<b>MiniGPT-v2</b> (Chen et al., 2023d) is an MM-LLM designed as a unified interface for diverse VL multi-task learning. To create a single model proficient in handling multiple VL tasks, identifiers are incorporated for each task during both training and inference. This facilitates clear task distinction, ultimately enhancing learning efficiency.	2467
2427		2468
2428		2469
2429		2470
2430		2471
2431		2472
2432		2473
2433	<b>CogVLM</b> (Wang et al., 2023b) is an open-source MM-LLM that bridges the gap between modalities via a trainable visual expert module within the attention and feedforward layers. This allows for a deep fusion of MM features without compromising performance on NLP downstream tasks.	2474
2434		2475
2435		2476
2436		2477
2437		2478
2438		2479
		2480
		2481
		2482
	<b>DRESS</b> (Chen et al., 2023i) introduces a method using natural language feedback to enhance alignment with human preferences. DRESS extends the conditional reinforcement learning algorithm to integrate non-differentiable natural language feedback, training the model to generate appropriate responses based on feedback.	2483
		2484
		2485
	<b>X-InstructBLIP</b> (Panagopoulou et al., 2023) introduces a cross-modal framework with instruction-aware representations, scalable enough to empower LLMs to handle diverse tasks across multiple modalities, including image/video, audio, and 3D. Notably, it achieves this without the need for modality-specific PT.	2486
		2487
		2488
		2489
		2490
		2491
		2492
	<b>CoDi-2</b> (Tang et al., 2023b) is a MM generation model excelling in modality-interleaved instruction following, in-context generation, and user-model interaction by multi-turn conversations. It enhances CoDi (Tang et al., 2023c) to process intricate modality-interleaved inputs and instructions, generating latent features autoregressively.	2493
		2494
		2495
		2496
		2497
		2498
		2499
	<b>VILA</b> (Lin et al., 2023) outperforms in vision tasks and shows remarkable reasoning ability while maintaining text-only capabilities. It achieves this by harnessing the full capabilities of LLM learning, using the interleaved attributes of image-text pairs, and implementing meticulous text data re-blending.	2500
		2501
		2502
		2503
		2504
		2505
		2506
		2507
		2508
		2509
		2510
		2511
		2512
		2513
		2514
		2515
		2516
		2517
		2518
		2519
		2520
		2521
		2522
		2523
		2524
		2525
		2526
		2527
		2528
		2529
		2530
		2531
		2532
		2533
		2534
		2535
		2536
		2537
		2538
		2539
		2540
		2541
		2542
		2543
		2544
		2545
		2546
		2547
		2548
		2549
		2550
		2551
		2552
		2553
		2554
		2555
		2556
		2557
		2558
		2559
		2560
		2561
		2562
		2563
		2564
		2565
		2566
		2567
		2568
		2569
		2570
		2571
		2572
		2573
		2574
		2575
		2576
		2577
		2578
		2579
		2580
		2581
		2582
		2583
		2584
		2585
		2586
		2587
		2588
		2589
		2590
		2591
		2592
		2593
		2594
		2595
		2596
		2597
		2598
		2599
		2600

Dataset Name	X Modality	#.X	#.T	#.X-T
ALIGN (Jia et al., 2021)	Image	1.8B	1.8B	1.8B
LTIP (Alayrac et al., 2022)	Image	312M	312M	312M
MS-COCO (Lin et al., 2014)	Image	124K	620K	620K
Visual Genome (Krishna et al., 2017)	Image	108K	4.5M	4.5M
CC3M (Sharma et al., 2018)	Image	3.3M	3.3M	3.3M
CC12M (Changpinyo et al., 2021)	Image	12.4M	12.4M	12.4M
SBU (Ordóñez et al., 2011)	Image	1M	1M	1M
LAION-5B (Schuhmann et al., 2022)	Image	5.9B	5.9B	5.9B
LAION-400M (Schuhmann et al., 2021)	Image	400M	400M	400M
LAION-en (Schuhmann et al., 2022)	Image	2.3B	2.3B	2.3B
LAION-zh (Schuhmann et al., 2022)	Image	142M	142M	142M
LAION-COCO (Schuhmann et al., 2022b)	Image	600M	600M	600M
Flickr30k (Young et al., 2014)	Image	31K	158K	158K
AI Challenger Captions (Wu et al., 2017)	Image	300K	1.5M	1.5M
COYO (Byeon et al., 2022)	Image	747M	747M	747M
Wukong (Gu et al., 2022)	Image	101M	101M	101M
COCO Caption (Chen et al., 2015)	Image	164K	1M	1M
WebLI (Chen et al., 2022b)	Image	10B	12B	12B
Episodic WebLI (Chen et al., 2023h)	Image	400M	400M	400M
CC595k (Liu et al., 2023e)	Image	595K	595K	595K
RefCOCO (Kazemzadeh et al., 2014)	Image	20K	142K	142K
RefCOCO+ (Yu et al., 2016)	Image	20K	142K	142K
Visual-7W (Zhu et al., 2016)	Image	47.3K	328K	328K
OCR-VQA (Mishra et al., 2019)	Image	207K	1M	1M
ST-VQA (Biten et al., 2022)	Image	23K	32K	32K
DocVQA (Mathew et al., 2021)	Image	12K	50K	50K
TextVQA (Singh et al., 2019)	Image	28.4K	45.3K	45.3K
DataComp (Gadre et al., 2023)	Image	1.4B	1.4B	1.4B
GQA (Hudson and Manning, 2019)	Image	113K	22M	22M
VGQA (Krishna et al., 2017)	Image	108K	1.7M	1.7M
VQA <sup>v2</sup> (Goyal et al., 2017)	Image	265K	1.4M	1.4M
DVQA (Kafle et al., 2018)	Image	300K	3.5M	3.5M
OK-VQA (Schwenk et al., 2022)	Image	14K	14K	14K
A-OKVQA (Schwenk et al., 2022)	Image	23.7K	24.9K	24.9K
Text Captions (Sidorov et al., 2020)	Image	28K	145K	145K
M3W (Interleaved) (Alayrac et al., 2022)	Image	185M	182GB	43.3M (Instances)
MMC4 (Interleaved) (Zhu et al., 2023c)	Image	571M	43B	101.2M (Instances)
Obelics (Interleaved) (Laurençon et al., 2023)	Image	353M	115M	141M (Instances)
MSRVTT (Xu et al., 2016)	Video	10K	200K	200K
WebVid (Bain et al., 2021)	Video	10M	10M	10M
VTP (Alayrac et al., 2022)	Video	27M	27M	27M
AISHELL-1 (Chen et al., 2023b)	Audio	-	-	128K
AISHELL-2 (Chen et al., 2023b)	Audio	-	-	1M
WaveCaps (Mei et al., 2023)	Audio	403K	403K	403K
VSDial-CN (Ni et al., 2024)	Image, Audio	120K (Image), 1.2M(Audio)	120K	1.2M

Table 3: The statistics for MM PT datasets. **#.X** represents the quantity of X, **#.T** represents the quantity of Text, and **#.X-T** represents the quantity of X-Text pairs, where X can be Image, Video, or Audio.

Dataset Name	Type	I→O	Source	Method	Multi-Turn	#I/V/A	#Dialog Turn	#Instance
MiniGPT-4's IT (Zhu et al., 2023a)	SFT	I→T	CC3M, CC12M	Auto.	✗	134M/-/-	1	5K
StableLLaVA (Li et al., 2023i)	SFT	I→T	SD (Rombach et al., 2022)	Auto.+Manu.	✗	126K/-/-	1	126K
LLaVA's IT (Zhang et al., 2023h)	SFT	I→T	MS-COCO	Auto.	✓	81K/-/-	2.29	150K
SVIT (Zhao et al., 2023a)	SFT	I→T	MS-COCO, Visual Genome	Auto.	✓	108K/-/-	5	3.2M
LLaVAR's IT (Zhang et al., 2023h)	SFT	I→T	MS-COCO, CC3M, LAION	LLaVA+Auto.	✓	20K/-/-	2.27	174K
ShareGPT4V's IT (Chen et al., 2023f)	SFT	I→T	LCS, COCO, SAM, TextCaps, WikiArt	Auto.+Manu.	✗	100K/-/-	-	-
DRESS's IT (Chen et al., 2023i)	SFT	I→T	LLaVA's IT, VLSafe	Auto.+Manu.	✓	193K/-/-	~4	-
VideoChat's IT (Li et al., 2023f)	SFT	V→T	WebVid	Auto.	✓	~8K/-	1.82	11K
Video-ChatGPT's IT (Maaz et al., 2023)	SFT	V→T	ActivityNet (Caba Heilbron et al., 2015)	Inherit	✓	~100K/-	1	100K
Video-LLaMA's IT (Zhang et al., 2023e)	SFT	I/V→T	MiniGPT-4, LLaVA, and VideoChat's IT	Auto.	✓	81K/8K/-	2.22	171K
InstructBLIP's IT (Dai et al., 2023)	SFT	I/V→T	Multiple (InstructBLIP's Figure 2)	Auto.	✗	-	-	~1.6M
X-InstructBLIP's IT (Panagopoulou et al., 2023)	SFT	I/V/A/3D→T	Multiple (X-InstructBLIP's Figure 4)	Auto.	✗	-	-	~1.8M
MIMIC-IT (Li et al., 2023a)	SFT	I/V→T	Multiple	Auto.	✗	8.1M/502K/-	1	2.8M
PandaGPT's IT (Su et al., 2023)	SFT	I→T	MiniGPT-4 and LLaVA's IT	Inherit	✓	81K/-/-	2.29	160K
MGVLD (Zhao et al., 2023b)	SFT	I+B→T	Multiple	Auto.+Manu.	✗	108K/-/-	-	108K
M <sup>3</sup> IT (Li et al., 2023h)	SFT	I/V/B→T	Multiple	Auto.+Manu.	✗	-/-/-	1	2.4M
LAMM (Yin et al., 2023b)	SFT	I+3D→T	Multiple	Auto.+Manu.	✓	91K/-/-	3.27	196K
BuboGPT's IT (Zhao et al., 2023d)	SFT	(I+A)/A→T	Clotho, VGGSS	Auto.	✗	5K/-/9K	-	9K
mPLUG-DocOwl's IT (Ye et al., 2023b)	SFT	I/Tab/Web→T	Multiple	Inherit	✗	-	-	-
T2M (Wu et al., 2023d)	SFT	T→I/V/A→T	WebVid, CC3M, AudioCap	Auto.	✗	4.9K/4.9K/4.9K	1	14.7K
MosIT (Wu et al., 2023d)	SFT	I+V+A+T→I+V+A+T	Youtube, Google, Flickr30k, Midjourney, etc.	Auto.+Manu.	✓	4K/4K/4K	4.8	5K
Osprey's IT (Yuan et al., 2023a)	SFT	I→T	MS-COCO, RefCOCO, RefCOCO+, LLaVA's IT etc. (fine-grained region-text dataset)	Auto.+Manu.	✓	-/-/-	~4	724K
LLaVA-RLHF (Sun et al., 2023b)	RLHF	I→T	Collected human preference	Manu.	✗	-/-/-	-	10K
DRESS's IT (Chen et al., 2023i)	RLHF	I→T	LLaVA's IT, VLSafe	Auto.+Manu.	✓	33K/-/-	~4	-
RLHF-V's IT (Yu et al., 2023b)	RLHF	I→T	Collected human preference	Manu.	✗	-/-/-	-	1.4K
VLFeedback (Li et al., 2023g)	RLHF	I→T	Responses generated by 12 MM-LLMs	Auto.	✗	-/-/-	-	80K
RTVLM (Li et al., 2024a)	RLHF	I→T	New question-image pairs based on publicly available images or originally diffusion-generated images (Gallegos et al., 2023)	Auto.+Manu.	✗	-/-/-	-	5K
VLGuard's IT (Zong et al., 2024)	RLHF	I→T	Source image data from various datasets	Auto.	✗	3K/-/-	-	3K
MMVG (Yan et al., 2024)	RLHF	I→T	MS-COCO	Manu.	✗	16K/-/-	-	16K

Table 4: The statistics for MM IT datasets. I→O: Input to Output Modalities, T: Text, I: Image, V: Video, A: Audio, B: Bounding box, 3D: Point Cloud, Tab: Table, and Web: Web page.