# ImageBind-LLM: Multi-modality Instruction Tuning

Jiaming Han[1,2†], Renrui Zhang[1,2†], Wenqi Shao[1†], Peng Gao[1‡†], Peng Xu[1†], Han Xiao[1†], Kaipeng Zhang[1], Chris Liu[1], Song Wen[1], Ziyu Guo[1], Xudong Lu[1,2], Shuai Ren[3], Yafei Wen[3], Xiaoxin Chen[3], Xiangyu Yue[2*], Hongsheng Li[2*], Yu Qiao[1*]

[1]Shanghai Artificial Intelligence Laboratory, Shanghai, 200030, China.
[2]CUHK MMLab, Hong Kong SAR, 999077, China.
[3]vivo AI Lab, Shenzhen, 518000, China.

*Corresponding author(s). E-mail(s): xyyue@ie.cuhk.edu.hk; hsli@ee.cuhk.edu.hk; qiaoyu@pjlab.org.cn;
Contributing authors: hanjiaming@pjlab.org.cn; zhangrenrui@pjlab.org.cn; shaowenqi@pjlab.org.cn; gaopeng@pjlab.org.cn; xupeng@pjlab.org.cn;
†Equal Contribution   ‡Project Leader

## Abstract

We present **ImageBind-LLM**, a multi-modality instruction tuning method of large language models (LLMs) via ImageBind. Existing works mainly focus on language and image instruction tuning, different from which, our ImageBind-LLM can respond to multi-modality conditions, including audio, 3D point clouds, video, and their embedding-space arithmetic by only image-text alignment training. During training, we adopt a learnable bind network to align the embedding space between LLaMA and ImageBind's image encoder. Then, the image features transformed by the bind network are added to word tokens of all layers in LLaMA, which progressively injects visual instructions via an attention-free and zero-initialized gating mechanism. Aided by the joint embedding of ImageBind, the simple image-text training enables our model to exhibit superior multi-modality instruction-following capabilities. During inference, the multi-modality inputs are fed into the corresponding ImageBind encoders, and processed by a proposed visual cache model for further cross-modal embedding enhancement. The training-free cache model retrieves from three million image features extracted by ImageBind, which effectively mitigates the training-inference modality discrepancy. Notably, with our approach, ImageBind-LLM can respond to instructions of diverse modalities and demonstrate significant language generation quality. Code is released at https://github.com/OpenGVLab/LLaMA-Adapter.

**Keywords:** Large Language Model, Multi-Modal Learning, Instruction Tuning

## 1 Introduction

Recently, we have witnessed substantial advancements in the instruction tuning of large language models (LLMs). With versatile intelligence and interactivity, ChatGPT [4] and GPT-4 [5] present general-purpose chatting systems following human instructions in language and images, which is yet unreplicable due to the closed-source restriction. Inspired by this, Alpaca [6], LLaMA-Adapter [7], and follow-up works [8–10] propose to fine-tune the publicly available LLaMA [11] into
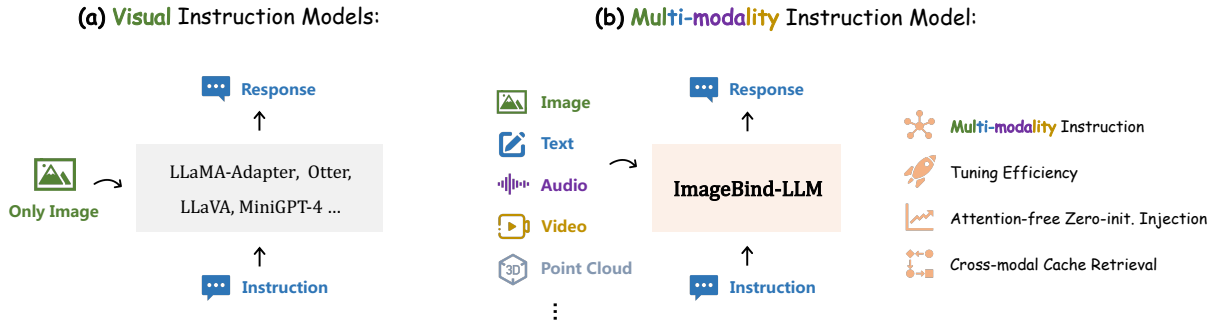
**Fig. 1 Visual Instruction Models vs. Our Multi-modality ImageBind-LLM.** Different from existing works [1–3] conditioned only on image modality, ImageBind-LLM conducts a general multi-modality instruction tuning for image, text, audio, video, and 3D.

language instruction models by self-constructed data. Further, to achieve image instruction tuning, LLaVA [2], LLaMA-Adapter [7], and others [3] incorporate visual understanding capabilities into LLMs for image-conditioned generation. Despite the effectiveness of existing instruction tuning approaches, how to develop an LLM for general multi-modality instructions, e.g., text, image, audio, 3D point clouds, and video, is still under-explored.

In this paper, we introduce a multi-modality instruction-following model, **ImageBind-LLM**, which efficiently fine-tunes LLaMA, guided by the joint embedding space in the pre-trained Image-Bind [12]. As compared in Figure 1, different from previous visual instruction models (a), our ImageBind-LLM (b) can respond to input instructions of multiple modalities besides images, indicating promising extensibility and generalization capacity. Specifically, thanks to the image-aligned multi-modality embedding space of ImageBind, we propose to only leverage the vision-language data for multi-modality instruction tuning. For an image-caption pair, we first utilize the frozen image encoder of ImageBind to extract the global image feature, and adopt a learnable bind network for embedding transformation. Then, the transformed image feature is added to the word tokens at all transformer layers in LLaMA, which provides visual conditions to generate the corresponding textual caption. Different from the zero-initialized attention in LLaMA-Adapter series [1, 7], our visual injection method is attention-free and simply weighted by a trainable zero-initialized gating factor. In such an efficient manner, the instruction cues of ImageBind's multi-modality

embeddings can be progressively injected into LLaMA as the training goes on, without disturbing the original language knowledge.

After the simple vision-language training, our ImageBind-LLM obtains the capability to follow instructions of various modalities, by applying ImageBind for modality-specific encoding, e.g., text, image, audio, and video. For instructions in 3D domains, we utilize the pre-trained 3D encoder in Point-Bind [13] to encode the input 3D point clouds. To alleviate the modality discrepancy of image training and text/audio/3D/video-conditioned generation, we further propose a training-free visual cache model for embedding enhancement during inference. The cache model contains millions of image features in the training datasets extracted by ImageBind, which improves text/audio/3D/video embeddings by retrieving similar visual features, referring to Tip-Adapter [14]. This contributes to higher-quality language responses to multi-modality instructions. In diverse scenarios, we evaluate the multi-modality instruction-following capabilities of ImageBind-LLM, and observe consistent superior performance.

Overall, our ImageBind-LLM exhibits four main characteristics as follows.

- **Multi-modality Instructions.** Different from previous language and image instruction models, ImageBind-LLM is tuned to respond to general multi-modality inputs, such as image, text, audio, 3D point clouds, video, and their embedding-space arithmetic encoded by Image-Bind and Point-Bind.

2

- **Tuning Efficiency.** During training, we freeze the image encoder of ImageBind, and fine-tune partial weights in LLaMA by parameter-efficient techniques, including LoRA [15] and bias-norm tuning [1, 16–19]. Besides, we only train the additional bind network and zero-initialized gating factors.
- **Attention-free Zero-initialized Injection.** Instead of incorporating new instruction cues by attention layers, we directly add the multi-modality conditions with all word tokens of LLaMA, and adopt a learnable gating mechanism for progressive knowledge injection, more simple and effective.
- **Cross-modality Cache Retrieval.** To alleviate the modality discrepancy of training (only image) and inference (multiple modalities), we introduce a visual cache model constructed by ImageBind-extracted image features, which conducts cross-modality retrieval for embedding enhancement.

## 2  Related Work

### 2.1  Visual Instruction Models.

Given the rapid development of language instruction-following capabilities [6, 8, 11], how to enable large language models (LLMs) to perform visual understanding has also gained significant attention. LLaMA-Adapter [7], for the first time, proposes to generate language responses conditioned on image inputs. It leverages a pre-trained encoder to extract image tokens, and incorporates them with LLaMA by parameter-efficient fine-tuning, which however can only tackle some naive visual question answering scenarios, i.e., ScienceQA [20]. For more general visual instruction-following circumstances, many efforts have been made to produce high-quality vision-language data for training by ChatGPT [4] or GPT-4 [21], such as LLaVA [22], MiniGPT-4 [3], and Otter [23]. They normally follow the architecture of BLIP-2 [24] with a more advanced Vicuna [9], or fine-tune the entire LLM with costly training resources. LLaMA-Adapter [7] develops a joint training strategy that only requires a combination of image-caption pairs and language instruction data, but still performs comparably to those with delicately constructed training data. VideoLLM [25] and Video-LLaMA [26] also connect video reasoning modules with LLMs to allow for video instruction-following powers with temporal information. Different from them, our ImageBind-LLM takes a step forward by tuning a multi-modality LLM conditioned on language questions with image, video, audio, and 3D point cloud input, allowing for widespread applications.

### 2.2  Multi-modality Alignment.

Bridging different modalities within a joint embedding space for cross-modality processing has emerged as a critical research area in both vision and language. CLIP [27], ALIGN [28], and Florence [29] utilize simple contrastive learning paradigms to align image and text pairs, contributing to promising zero-shot generalization performance. Flamingo [30], BLIP-2 [24], and MAGIC [31] adopt intermediate networks to connect pre-trained vision and language encoders. AudioCLIP [32] and PointCLIP [33] respectively extend the embedding space of CLIP to other modalities, such as audio and 3D point clouds. Recently, ImageBind [12] is proposed to share a single latent space with various modalities, including image, video, text, and audio. Inspired by ImageBind, Point-Bind [13] learns to blend 3D point cloud modalities into ImageBind, and achieves favorable 3D zero-shot accuracy. In this paper, we focus on aligning the shared embedding space in ImageBind/Point-Bind with LLaMA for multi-modality instruction-following capacity. PandaGPT [34] also aims to tune a multi-modality LLM based on ImageBind, which cannot support 3D point clouds as input, and utilizes a stronger LLM, Vicuna [9], as the pre-trained language model. In contrast, our ImageBind-LLM is still based on LLaMA [11] and introduces unique attention-free zero-initialized injection with cross-modality cache retrieval for better multi-modality reasoning.

## 3  Method

In Section 3.1, we first briefly revisit some prior works as a preliminary, including ImageBind, cache models, and LLaMA-Adapter. Then, in Section 3.2, we introduce the details of our proposed multi-modality instruction tuning and cache-enhanced inference in ImageBind-LLM.

## 3.1 A Revisit of Prior Works

### 3.1.1 ImageBind

With a single joint embedding space, Image-Bind [12] proposes to connect five different modalities, i.e., text, audio, depth, thermal, and Inertial Measurement Unit (IMU), all by image-paired data. Following CLIP [27], the pre-training of ImageBind adopts a contrastive loss, which clusters image features with other paired modalities, and pushes away unpaired ones in the embedding space. Self-supervised by large-scale image-paired data, ImageBind learns to encode different modalities into aligned feature embeddings, which obtains emergent cross-modal zero-shot capabilities. Then, ImageBind can be utilized to extend existing vision-language models to incorporate new modalities, such as text-to-audio/video retrieval, audio-to-image generation, and audio-referred object detection. Inspired by this image-centric property, our approach only conducts vision-language training to align the joint embedding space of ImageBind with LLaMA [11], achieving efficient multi-modality instruction tuning.

### 3.1.2 LLaMA-Adapter

As a novel parameter-efficient fine-tuning method, LLaMA-Adapter [7] transforms LLaMA into a language instruction model by only 1.2M parameters within 1 hour, which exhibits comparable performance to the fully fine-tuned Alpaca [6]. On top of this, LLaMA-Adapter [7] is also proposed to attain superior visual instruction-following capacity. It adopts a joint training paradigm for image-text and language-only instruction data, and still features tuning efficiency by updating partial parameters (14M) in LLaMA. One of the core innovations of LLaMA-Adapter series is the zero-initialized attention mechanism. They encode vision instruction signals as tokens, and concatenate them with the word tokens in LLaMA as prefixes. Within every attention layer, a learnable gating factor is utilized to adaptively control how much information the new instruction knowledge is incorporated into LLMs. Our ImageBind-LLM also adopts a zero-gated injection strategy for multi-modality instructions, but in a more simple and effective attention-free manner.

### 3.1.3 Cache Models

Without any training, a cache model can be utilized to store the features and labels of a training set, organizing them as a key-value database. During inference, the test sample serves as a query to retrieve from the keys and aggregate informative values via the key-query similarity. Starting from the conventional $k$ Nearest Neighbors algorithm ($k$-NN), cache models have been widely adopted to assist deep neural networks in language [35], 2D vision [36], and 3D point clouds [37]. Tip-Adapter [14] and its follow-up works [37–39] propose to store the CLIP-extracted image features of the given few-shot data, and regard the cache model as a non-parametric adapter for downstream tasks. Similarly, we cache the ImageBind-extracted 1 million image features as both keys and values, which enhances the multi-modality embeddings in inference time.

## 3.2 ImageBind-LLM

To obtain a multi-modality instruction model, we propose ImageBind-LLM, which includes two training stages: vision-language pre-training on image-caption data (Section 3.2.1) and multi-modality instruction tuning on visual instruction data (Section 3.2.2). Besides, we also propose cross-modality cache retrieval for enhanced inference (Section 3.2.3). The overall training paradigm of ImageBind-LLM is shown in Figure 2.

### 3.2.1 Vision-Language Pre-training

Given the modality-bound property of Image-Bind [12], we only fine-tune LLaMA [11] to generate language responses conditioned on ImageBind-encoded images, after which, the model can inherently understand instructions of other modalities via the respective ImageBind encoders. Therefore, we propose to only leverage vision-language data for tuning a multi-modality instruction model. Following LLaMA-Adapter [7], we adopt a two-stage training pipeline for ImageBind-LLM: first utilizing large-scale image-caption data [40–42] to learn the image-conditioned response capacity, then leveraging instruction-following data [3, 22] to preserve the long-sentence generation quality. The overall training paradigm of ImageBind-LLM is shown in Figure 2. For a given image-caption pair, we first adopt the frozen image encoder
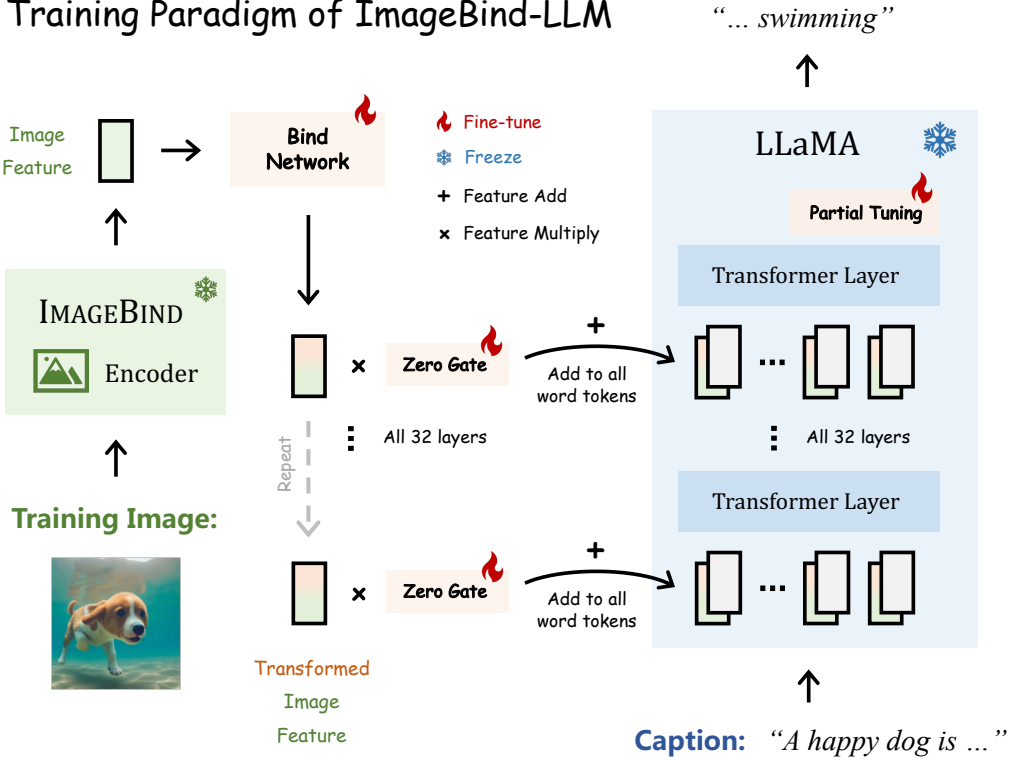
**Fig. 2 Training Paradigm of ImageBind-LLM.** Considering the joint embedding space in imageBind [12], we only utilize image-text datasets for multi-modality instruction tuning of LLaMA [11]. We introduce a bind network for embedding alignment, and an attention-free zero-initialized mechanism for visual knowledge injection.

of ImageBind to extract the global visual feature. Then, we transform the visual feature with a learnable bind network, and add it to every word token in LLaMA. In an attention-free zero-initialized manner, LLaMA is injected by image condition and generates the given image caption.

**Bind Network.**

In Figure 3, we present the details of the bind network, which aims to align the embedding space between ImageBind and LLaMA. Specifically, we denote the $C_I$-dimensional global image feature encoded by ImageBind as $F_I \in \mathbb{R}^{1 \times C_I}$. In the bind network, we first adopt a linear projection layer with a weight matrix $w_0 \in \mathbb{R}^{C_I \times C}$, formulated as $F_I^0 = F_I w_0 \in \mathbb{R}^{1 \times C}$, where $C$ denotes the feature dimension of LLaMA. Inspired by the Feed-Forward Network (FFN) in LLaMA, we then cascade three projection blocks with RMSNorm [43], SiLU activation functions [44], and residual connections [45]. For the $(i + 1)$-th

block with $F_I^i$ as input, we formulate the calculation of $F_I^{i+1}$ as (the normalization is omitted for simplicity)

$$F_I^{i+1} = F_I^i + (F_I^i w_2 \cdot \text{SiLU}(F_I^i w_1))w_3, \quad 0 \leq i < 3 \tag{1}$$

where $w_1, w_2 \in \mathbb{R}^{C \times C_h}$ and $w_3 \in \mathbb{R}^{C_h \times C}$, with $C_h$ denoting the hidden dimension. After the bind network, we obtain the transformed image feature, $T_I \in \mathbb{R}^{1 \times C}$, which learns to align the embedding space from ImageBind to LLaMA.

**Attention-free Zero-initialized Injection.**

With the encoded image feature $T_I$, existing visual instruction methods, e.g., LLaMA-Adapter [7], LLaVA [22], and MiniGPT-4 [3], concatenate it as the prefix to the word token sequence $\{T_W^j\}_{j=1}^N$ in LLaMA, where $N$ denotes the sequence length. Then, they leverage self-attention mechanisms in LLaMA's transformer blocks for visual knowledge incorporation from $T_I$ to $\{T_W^j\}_{j=1}^N$. However, such
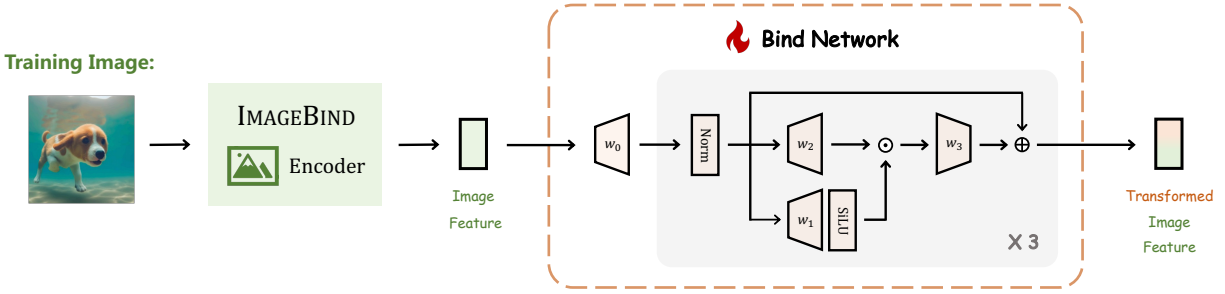
**Fig. 3 Details of the Bind Network.** Referring to the Feed-Forward Network (FFN) in LLaMA [11], we adopt cascaded blocks of RMSNorm [43], SiLU activation functions [44], and residual connections [45]. This aims to align the image feature from ImageBind [12] with LLaMA's word embeddings.

an attention-based approach not only causes extra computation budget, but also increases the training difficulty. In our ImageBind-LLM, we adopt a simpler and more effective method by attention-free zero-initialized injection. We directly add the image feature $T_I$ with every word token at all transformer layers of LLaMA, which explicitly fuses the visual conditions (and multi-modality inputs during inference) with the language knowledge in LLM. In addition, to adaptively control the level of integration, we utilize a learnable gating factor initialized by zero, denoted as $g_{zero}$. For any word token $T_W^j$ in LLaMA, we formulate the visual injection as

$$T^j = T_I \cdot g_{zero} + T_W^j. \tag{2}$$

Similar to the zero-initialized attention in LLaMA-Adapter [7], this gating factor can progressively increase during training, and inject more visual semantics into LLaMA, contributing to stable learning in the early training stage.

### 3.2.2 Multi-modality Instruction Tuning

Since we have connected ImageBind and LLaMA with a bind network via large-scale image-text pre-training, ImageBind-LLM can understand multimodal inputs (audio, video, and 3D point clouds), and generate language response conditioned on multi-modality inputs. However, unlike LLaVA [22] and MiniGPT-4 [3] that directly utilize a well-trained language instruction model Vicuna [9] as the base LLM, we instead adopt a non-instruction model LLaMA. Therefore, in the second training stage, we partially tune the

parameters in LLaMA to equip it with instruction-following ability, while keep the multi-modality encoders of ImageBind and the bind network frozen.
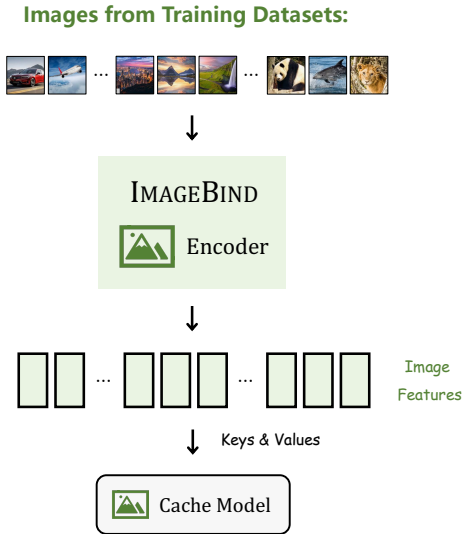
#### *Parameter-efficient Fine-tuning.*

Our second-stage training data is a mixture of language instruction data [10, 46] and visual instruction data [22]. As language instruction data contains no paired images, we input a fake image (filled with zero) as the placeholder during training. To maintain the multi-modality understanding capacity of the first-stage training, we only fine-tune LLaMA with parameter-efficient methods including Low-Rank Adaptation (LoRA) [15] and bias-norm tuning [1, 16–19]. Specifically, we add a low-rank layer for each linear layer in the transformer, where the rank is set to 16 by default. We also unfreeze all the normalization layers and add a learnable bias term to the linear layers. The parameters of all other modules are frozen during training.

#### *High-Quality Instruction Tuning.*

Although the fine-tuned ImageBind-LLM can generate instruction-following responses, we notice that it occasionally fantasizes about objects that don't exist in the input modality. Therefore, we introduce additional instruction tuning stage using high-quality instruction data from MiniGPT-4 [3]. Different from the visual instruction data generated by ChatGPT/GPT4,

**(a) Cache Model Construction**

**Images from Training Datasets:**

IMAGEBIND Encoder

Image Features

Keys & Values

Cache Model

**(b) Cache-Enhanced Inference**

**Multi-modality Input:**

IMAGEBIND Multi-modality Encoders

Naive Multi-modality Features

Cache Model

Query

Retrieve top-*k*

LLaMA

Bind Network

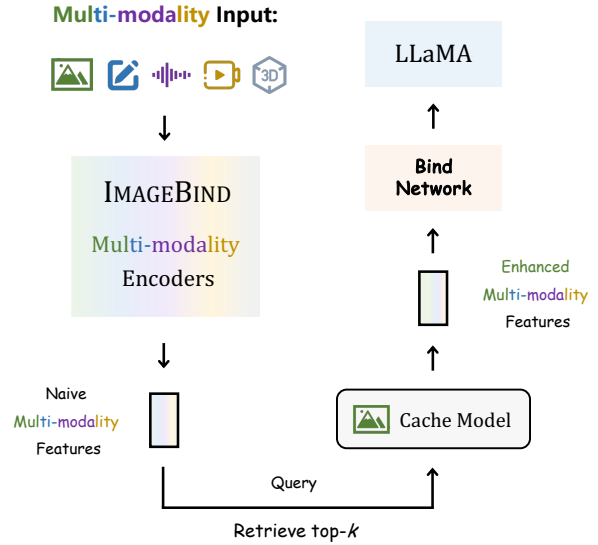Enhanced Multi-modality Features

**Fig. 4 Cache Retrieval for Inference.** To mitigate the training-inference discrepancy, we construct a training-free visual cached model of ImageBind-encoded image features (b). Then, during inference, we enhance the multi-modality embeddings by retrieving top-*k* similar visual features in the cache model.

MiniGPT-4 manually collects 3.5K image description data for high-quality vision-language alignment. Here we also adopt the 3.5K image description data for further instruction tuning, which only takes a few minutes.

### 3.2.3 Cache Retrieval for Inference

After fine-tuning by visual instruction data, ImageBind-LLM can inherently derive the multi-modality instruction-following capacity. Then, besides the naive inference mode (Figure 4 (a)), we further propose to construct a visual cache model by ImageBind for multi-modality embedding enhancement (Figure 4 (b)).

*Naive Multi-modality Inference.*

Via the bind network, the image features from ImageBind can be well aligned with LLaMA's word tokens. Given the joint multi-modality embedding space of ImageBind, our ImageBind-LLM spontaneously obtains the understanding capability for input instructions with various modalities, such as image, text, audio, and video, only if we apply their corresponding encoders from ImageBind before the bind network. For 3D point clouds, we can also utilize the pre-trained

3D encoder of Point-Bind [13] for global feature extraction, which shares the same embedding space with ImageBind.

*Cache-enhanced Inference.*

Despite the effectiveness of the aforementioned naive mode, there exists modality discrepancy in ImageBind-LLM between training and inference. Namely, we adopt image encoder of ImageBind for training, but switch to other encoders for inference, which slightly disturbs the tuned bind network and LLaMA. Therefore, we construct a training-free cache model of image features to enhance the multi-modality embeddings during inference. As shown in Figure 4 (a), we utilize ImageBind to encode a subset of the vision-language training data, and store them as both keys and values in the cache model. For an input multi-modality instruction in Figure 4 (b), we regard its ImageBind-encoded feature as the query, $F_M \in \mathbb{R}^{1 \times C_I}$, and retrieve the top-*k* similar visual keys from the cache model, denoted as $F_{key} \in \mathbb{R}^{k \times C_I}$. We formulate the top-*k* cosine similarity as

$$S_{topk} = F_M F_{key}^T \quad \in \mathbb{R}^{1 \times k}, \tag{3}$$

where we suppose $F_M$ and $F_{key}$ have been L2-normalized. Then, according to $S_{topk}$, we aggregate the corresponding cached values, $F_{value} \in \mathbb{R}^{k \times C_I}$ (top-$k$ similar image features), and add the result to the original feature $F_M$ via a residual connection, formulated as

$$F_M^e = \alpha \cdot S_{topk} F_{value} + (1 - \alpha) \cdot F_M, \quad (4)$$

where $\alpha$ serves as a balance factor. Aided by the cache model, the enhanced feature $F_M^e$ is adaptively incorporated with similar visual semantics from the cache model. This boosts the representation quality of other modalities, and mitigates their semantic gap to the images used for training. After this, $F_M^e$ is fed into the bind network for feature transformation and LLaMA for response generation.

## 3.3 Advanced Applications

Besides the superior multi-modality instruction-following capabilities, our ImageBind-LLM can also be extended to a wide range of advanced applications with simple modifications.

### 3.3.1 Bilingual Instruction Tuning

In addition to English instructions, ImageBind-LLM can be easily upgraded to a bilingual instruction-following model, e.g., English and Chinese. More specifically, we replace the basic LLM from LLaMA to a bilingual LLM, ChineseLLaMA[1] and add 52K Chinese instruction data from GPT4LLM [10] for joint instruction tuning. Although we do not have direct Chinese visual instruction data for the first vision-language training stage, we observe that our bilingual ImageBind-LLM implicitly learns the alignment between Chinese, English and multi-modality inputs, and can well follow Chinese instructions conditioned on other modality inputs.

### 3.3.2 Any-to-any Generation

Currently, most multi-modality instruction models are limited to generating only textual responses, lacking the ability to respond with other modal outputs, e.g., image, audio, and point clouds. Since ImageBind is an extension of CLIP [27], we can append CLIP-conditioned generative models after ImageBind's encoders, such as Stable Diffusion [47], Make-An-Audio [48], and CLIP-Forge [49], respectively for image, audio, and point cloud generation. Instead of directly inputting ImageBind features into these generative models, we adopt cache-enhanced generation to mitigate the modality discrepancy, similar to the approach in Cache-enhanced Inference (Section 3.2.3). In this way, we can achieve instruction models with any-to-any generation system, i.e., responding to multi-modality instructions by multi-modality responses. as an example, our ImageBind-LLM can generate both textual and image responses for multi-modality inputs (*e.g.*, image, audio and point clouds).

### 3.3.3 Integration with Object Detection

Visual instruction models can answer questions based on the global content of input images. However, they cannot associate the text response with regional objects in the image, which is important for fine-grained tasks such as visual reasoning and grounding. We provide a solution to connect ImageBind-LLM with object detectors [50]. For a response generated by ImageBind-LLM, we use traditional noun parsers [51] or ChatGPT [4] to extract nouns in the response. Then we feed the input image and parsed nouns into object detectors to get object detection results. Generally, the traditional noun parser is enough for parsing meaningful nouns, but it cannot handle nouns with complex modifiers, such as "a running black dog". Therefore, we will also ask ChatGPT to extract complex nouns in the response.

### 3.3.4 ImageBind-LLM as Chatbot

ImageBind-LLM was originally designed as a single-turn multi-modality instruction model. We turn ImageBind-LLM into a multi-turn chatbot by training it on multi-turn conversation data, including language conversation data from ShareGPT [46] and visual conversation data from LLaVA [2]. By this, ImageBind-LLM can be used as a multi-turn chat model to answer open-ended questions on multi-modality inputs.

---

[1] https://github.com/OpenLMLab/OpenChineseLLaMA

### 3.3.5 ImageBind-LLM for API Control

In addition to its primary multimodal instruction-following capacity, ImageBind-LLM also exhibits the potential to invoke diverse API controls for multi-modality tool usage. To achieve this, we leverage the tool-related instruction dataset introduced in GPT4Tools [52] to empower ImageBind-LLM with the ability to effectively utilize various tools. By training ImageBind-LLM on the GPT4Tools dataset using our proposed training paradigm, we observe its impressive proficiency in calling different APIs, enabling it to accomplish a wide range of tasks, even when encountering previously unseen tools. This performance in API control highlights the potential of ImageBind-LLM as a versatile visual assistant capable of solving diverse real-world problems.

## 4 Experiment

### 4.1 Training Details

#### 4.1.1 Datasets

We train ImageBind-LLM on a collection of open-sourced image-text pair data, language-only and visual instruction data.

**Image-Text Pair Data.** Our ImageBind-LLM is pre-trained on the concatenation of open-sourced image-text pair data, including COCO [53], CC3M [41], CC12M [42], SBU [54], LAION-2B [40], COYO [55] and MMC4 [56]. Note that MMC4-Core [56] is a billion-scale corpus of images interleaved with text. We extract 20M high-quality image-text pairs from MMC4-Core according to the provided clip alignment score. For LAION-2B [40] dataset, we also extract 100M high-quality image-text pairs based on their CLIP alignment scores. The concatenation of all open-sourced image-text pairs result into 940M image-text pair data. Unlike BLIP [57] which designs an effective data cleaning pipeline, our image-text pairs are much noisy. However, we empirically observe strong image understanding and factual ability of ImageBind-LLM when pre-trained with this dataset. In the future, we will explore advanced approaches for data cleaning and deduplication.

**Instruction Tuning Datasets.** Our instruction tuning data includes language instruction data Alpaca [6], GPT4LLM [10] and ShareGPT [46], visual instruction data LLaVA [22] and MiniGPT4 [3]. For language instruction data, Alpaca contains 52K single-turn instruction data collected from GPT3.5; GPT4LLM is a GPT4 version of Alpaca with higher quality; ShareGPT is a collection of user-shared conversations with ChatGPT/GPT4. For visual instruction data, LLaVA adopts GPT4 to transform image captions or object detection annotations into 150K visual instruction data; MiniGPT4 curates a high-quality image description dataset with 3.5K examples. Note that we will convert multi-round conversation data into single turn data for instruction tuning.

#### 4.1.2 Implementation Details

For cache-enhanced inference, we use the FAISS library [58] to build our retrieval system, and the Autofaiss library[2] to find the optimal hyper-parameters for the index. By default, all images from CC3M [41] is used to build the cache model. We pre-train the model on 32 A100 GPUs for 3 epochs. The total batch size and learning rate is set to 1024 and 4e-4, respectively. We fine-tune the model on 8 A100 GPUs for 4 epochs The warmup epoch, total batch size, learning rate is set to 1, 32 and 1.25e-4.

### 4.2 Quantitative Evaluation on Traditional Tasks

In this section, we conducted quantitative evaluations of ImageBind-LLM on 27 datasets using a zero-shot approach. Our quantitative evaluation encompassed five specific tasks: Optical Character Recognition (OCR), Key Information Extraction (KIE), Image Captioning, Visual Question Answering (VQA), and Knowledge-Grounded Image Description (KGID). Notably, all these tasks are evaluated following a VQA-style approach. The comparisons of ImageBind-LLM with other well-known Vision-Language Models (VLMs) such as BLIP2 [24], InstructBLIP [60], LLaVA [22], LLaMA-Adapter (LA) [7], and multi-modality LLM model PandaGPT [34] are presented in Table 1 and Table 2.

---

[2]https://github.com/criteo/autofaiss

**Table 1 Zero-shot Performance on OCR, KIE, and Image Captioning Tasks.** Evaluation metrics include word accuracy for OCR datasets, entity-level F1 score for KIE datasets, and CIDEr score for image captioning datasets. ImageBind-LLM-D: ImageBind-LLM trained on multi-turn conversation data (Sec. 3.3.4).

| | Model | BLIP2 | InstructBLIP | LA | LLaVA | PandaGPT | ImageBind-LLM | ImageBind-LLM-D |
|---|---|---|---|---|---|---|---|---|
| | #Token | 32 | 32 | 10 | 257 | 1 | 1 | 1 |
| OCR | IIIT5K | 80.17 | 83.90 | 36.30 | 31.57 | 5.27 | 13.9 | 13.87 |
| | IC13 | 81.13 | 82.08 | 20.87 | 16.39 | 4.60 | 7.43 | 7.19 |
| | IC15 | 66.68 | 73.57 | 29.40 | 26.58 | 4.57 | 11.94 | 11.36 |
| | Total-Text | 68.31 | 71.51 | 30.93 | 24.51 | 4.06 | 10.79 | 10.11 |
| | CUTE80 | 85.07 | 86.11 | 35.76 | 36.46 | 6.60 | 20.14 | 20.83 |
| | SVT | 85.78 | 86.86 | 20.40 | 18.55 | 3.40 | 8.35 | 7.11 |
| | SVTP | 77.34 | 80.93 | 31.01 | 27.44 | 4.96 | 10.39 | 10.08 |
| | COCO-Text | 53.62 | 58.25 | 20.94 | 18.05 | 2.67 | 5.59 | 5.12 |
| | WordArt | 73.66 | 75.12 | 38.98 | 35.87 | 7.81 | 21.24 | 20.58 |
| | CTW | 67.43 | 68.58 | 18.13 | 16.73 | 2.74 | 7.12 | 7.38 |
| | HOST | 57.28 | 61.22 | 16.60 | 15.94 | 3.97 | 7.53 | 7.82 |
| | WOST | 68.83 | 73.26 | 21.73 | 20.49 | 4.01 | 8.73 | 8.57 |
| KIE | SROIE | 0.08 | 0.09 | 0.02 | 0.01 | 0.01 | 0.01 | 0.01 |
| | FUNSD | 1.02 | 1.03 | 2.16 | 1.93 | 2.06 | 2.00 | 2.01 |
| Caption | NoCaps | 48.58 | 46.33 | 41.66 | 33.09 | 29.65 | 30.43 | 29.64 |
| | Flickr-30k | 46.48 | 50.45 | 30.49 | 27.65 | 23.02 | 23.04 | 23.49 |

**Table 2 Zero-shot Performance on VQA, KGID, and VE Tasks.** For VQA and KGID tasks, Mean Reciprocal Rank (MRR) is used for the Visdial, while top-1 accuracy is employed for the remaining tasks.

| | Model | BLIP2 | InstructBLIP | LA | LLaVA | PandaGPT | ImageBind-LLM | ImageBind-LLM-D |
|---|---|---|---|---|---|---|---|---|
| | #Token | 32 | 32 | 10 | 257 | 1 | 1 | 1 |
| VQA | DocVQA | 4.75 | 5.89 | 8.13 | 6.26 | 3.42 | 4.04 | 4.08 |
| | TextVQA | 31.98 | 39.60 | 43.76 | 38.92 | 16.42 | 23.98 | 23.98 |
| | STVQA | 20.98 | 28.30 | 32.33 | 28.40 | 11.23 | 15.55 | 14.75 |
| | OCR-VQA | 38.85 | 60.20 | 38.12 | 23.40 | 22.39 | 23.24 | 22.31 |
| | OKVQA | 44.93 | 60.52 | 55.93 | 54.36 | 50.85 | 51.66 | 51.70 |
| | GQA | 45.53 | 49.96 | 43.93 | 41.30 | 41.56 | 41.23 | 41.12 |
| | Visdial | 10.73 | 45.20 | 12.92 | 14.66 | 90.80 | 12.66 | 12.91 |
| | IconQA | 62.82 | 56.25 | 41.83 | 42.95 | 46.04 | 37.97 | 41.81 |
| | VSR | 63.63 | 41.28 | 50.63 | 51.24 | 46.75 | 49.37 | 49.78 |
| KGID | ScienceQA IMG | 60.73 | 46.26 | 54.19 | 49.33 | 52.80 | 55.83 | 51.41 |
| | VizWiz | 65.44 | 65.31 | 62.07 | 62.42 | 46.95 | 51.90 | 51.28 |

## 4.2.1 Experimental Settings

**OCR Tasks.** We evaluate ImageBind-LLM on 12 representative OCR datasets, including IIIT5K [61], ICDAR 2013(IC13) [62], ICDAR 2015 (IC15) [63], Total-Text [64], CUTE80 [65], Street View Text (SVT) [66], SVTP-Perspective (SVTP) [67], COCO-Text [68], WordArt [69], SCUT-CTW1500 (CTW) [70], Heavily Occluded Scene Text (HOST) [71], Weakly Occluded Scene Text (WOST) [71]. These datasets encompass a diverse collection of images containing textual information, enabling a comprehensive comparison between models. The evaluation of model performance was based on top-1 accuracy, using the prompt "What is written in the image?"

**KIE Tasks.** We evaluate ImageBind-LLM on 2 KIE benchmarks, including SROIE [72] and FUNSD citefunsd. These benchmarks encompass a diverse range of document types, including receipts and forms, which necessitate the extraction of specific information. The evaluation of models involved using entity-level F1 scores. To further enhance the evaluation process, we employed prompts tailored to the specific information that the model was required to extract. For instance, in the case of the SROIE benchmark, prompts such as "What is the name of the

**Table 3 Perception Performance Comparison on MME [59] benchmark.** The full score for the overall perception tasks is 2000, while for the 10 subtasks is 200.

| Model | MiniGPT-4 | Otter | LLaMA-Adapter | LLaVA | PandaGPT | ImageBind-LLM |
|---|---|---|---|---|---|---|
| #Token | 32 | 64 | 10 | 257 | 1 | 1 |
| Existence | 115.00 | 48.33 | 120.00 | 50.00 | 70.00 | 128.33 |
| Count | 123.33 | 50.00 | 50.00 | 50.00 | 50.00 | 60.00 |
| Position | 81.67 | 50.00 | 48.33 | 50.00 | 50.00 | 46.67 |
| Color | 110.00 | 55.00 | 75.00 | 55.00 | 50.00 | 73.33 |
| Poster | 55.78 | 44.90 | 99.66 | 50.00 | 76.53 | 64.97 |
| Celerity | 65.29 | 50.00 | 86.18 | 48.82 | 57.06 | 76.47 |
| Scene | 95.75 | 44.25 | 148.50 | 50.00 | 118.00 | 113.25 |
| Landmark | 69.00 | 49.50 | 150.25 | 50.00 | 69.75 | 62.00 |
| Artwork | 55.75 | 41.75 | 69.75 | 49.00 | 51.25 | 70.75 |
| OCR | 95.00 | 50.00 | 125.00 | 50.00 | 50.00 | 80.00 |
| Perception | 866.58 | 483.73 | 972.67 | 502.82 | 642.59 | 775.77 |

**Table 4 Cognition Performance Comparison on MME [59] benchmark.** The full score for the overall perception tasks is 800, while for the 4 subtasks is 200.

| Model | MiniGPT-4 | Otter | LLaMA-Adapter | LLaVA | PandaGPT | ImageBind-LLM |
|---|---|---|---|---|---|---|
| #Token | 32 | 64 | 10 | 257 | 1 | 1 |
| Commonsense Reasoning | 72.14 | 38.57 | 81.43 | 57.14 | 73.57 | 48.57 |
| Numerical Calculation | 55.00 | 20.00 | 62.50 | 50.00 | 50.00 | 55.00 |
| Text Translation | 55.00 | 27.50 | 50.00 | 57.50 | 57.50 | 50.00 |
| Code Reasoning | 110.00 | 50.00 | 55.00 | 50.00 | 47.50 | 60.00 |
| Cognition | 292.14 | 136.07 | 248.93 | 214.64 | 228.57 | 213.57 |

company that issued this invoice?" were used to extract company information, while prompts like "Where was this invoice issued?" were employed to extract address information.

**VQA Tasks.** We employ 9 benchmarks in the VQA task, namely DocVQA [73], TextVQA [74], STVQA [75], OCR-VQA [76], OKVQA [77], GQA [78], IconQA [79], Visual Spatial Reasoning (VSR) [80], and Visual Dialog (Visdial) [81]. These benchmarks encompass a diverse collection of question-image pairs that cover a wide range of topics. The task requires models not only to comprehend the visual content but also to understand and reason about the questions presented. For specific evaluation purposes, we utilize the Mean Reciprocal Rank (MRR) metric for Visdial and top-1 accuracy for the remaining datasets. These metrics provide valuable insights into the model's proficiency in accurately answering questions across the various VQA benchmarks.

**KGID tasks.** The KGID task aims to assess the model's ability to produce descriptive and precise image captions by incorporating external knowledge. To evaluate performance in this task, we utilize the ScienceQA [20] and VizWiz [82] benchmarks, which include images accompanied by textual descriptions and knowledge-based information. It is worth mentioning that, for ScienceQA, we specifically consider only those samples that contain images.

### 4.2.2 Analysis

Table 1 and Table 2 clearly demonstrate the exceptional zero-shot performance of ImageBind-LLM across all evaluated tasks. When it comes to OCR, Image Captioning, and KGID, ImageBind-LLM achieved competitive performance compared with other VLMs and outperformed PandaGPT, thus showcasing the effectiveness of ImageBind-LLM's modality alignment strategy. Furthermore, ImageBind-LLM also delivered an impressive performance on KIE and VQA datasets.

Further investigating the reason behind ImageBind-LLM's relatively better performance

than PandaGPT, we delve into the implementation details of ImageBind-LLM and PandaGPT. Firstly, we observe a significant disparity in ImageBind-LLM and PandaGPT's utilization of the ImageBind extracted feature. PandaGPT employs a single linear projection layer for processing the ImageBind extracted feature, whereas ImageBind-LLM employs a bind network, which potentially facilitates better alignment between language and modalities through ImageBind. Another distinction lies in their choice of LLM model, with PandaGPT utilizing Vicuna and ImageBind-LLM employing LLaMA. Notably, Vicuna, being tuned based on LLaMA and possessing a higher Elo rating as indicated in [9], potentially enhances PandaGPT's language comprehension and response generation capabilities.

Then for why both ImageBind-LLM and PandaGPT have a poor OCR ability compared to other VLMs, we discovered that both of them employ only one token for the modality feature, while the other VLMs utilize at least ten tokens for capturing visual information. This disparity may allow other VLM models to better comprehend the visual information depicted in the images.

These results not only highlight the remarkable zero-shot performance of ImageBind-LLM in various vision and language tasks but also underscore its ability to comprehend and generate accurate responses in diverse scenarios. Moreover, the model's adeptness in multi-modality understanding further demonstrates its potential as a robust and versatile solution for real-world applications.

## 4.3 Quantitative Evaluation on MME Benchmark

### 4.3.1 Experimental Settings

In contrast to traditional multi-modality tasks, we also evaluate our ImageBind-LLM on a newly proposed benchmark, MME [59], which is specially deigned for the recent VLMs. MME benchmark systematically measures two multi-modality capabilities of existing methods: perception and cognition. The former with 10 subtasks refers to recognizing specific objects in images, while the latter with 4 subtasks is more challenging for deducing complex answers from visual information. For each test image, MME adopts an instruction of a question and a description "Please

answer yes or no", which prompts LLMs to answer "yes" or "no". Such a concise instruction-answer evaluation allows for fair comparison of LLMs without the impact of prompt engineering.

### 4.3.2 Analysis

In Table 3 and 4, we respectively show the performance comparison of different VLMs on MME's perception and cognition tasks, including MiniGPT-4 [3], Otter [23], LLaMA-Adapter [7], LLaVA [22], and PanadaGPT [34]. As shown, MiniGPT-4 can achieve the best scores since it is trained upon a pre-trained BLIP-2 [24]. Otter and PandaGPT are developed based on Open-Flamingo [83] and Vicuna [9], which endow them with well-initialized language processing abilities. Instead, similar to LLaMA-Adapter, our ImageBind-LLM is fine-tuned on the original LLaMA model, and still performs competitively to others. Especially on 'Existence' and 'Artwork', ImageBind-LLM outperforms the second-best methods by +8.33 and +1.00 scores, respectively. Overall, our approach is more expert at the 'Perception' tasks, ranking the third place and surpassing another multi-modality model, PandaGPT, by +133.18 score. As analyzed above in Section 3.2, we believe our performance can be further improved if using more multi-modality tokens fed into LLMs.

## 4.4 Qualitative Analysis

In this section, we will give qualitative examples and analysis to help understand how ImageBind-LLM works, and where its multi-modal instruction capabilities come from.

### 4.4.1 Multi-modality Understanding

**Multi-modality to Text Alignment.** The vision-language pre-training stage is essential for incorporating multi-modal information into LLMs. In Fig. 5, we give some multi-modality captioning results using the pre-trained ImageBind-LLM. As we can see, ImageBind-LLM can generate modality-dependent outputs for image, audio, video and point cloud. Since ImageBind-LLM is pre-trained with image-text pairs, it can give a short and accurate description of the image. Thanks to the binding property of ImageBind,

| Titanic | Abandoned Castle | Sleeping Fox | The City of the Future |

| Birds | Dog | Yoga in the Sun - 2018 | Airplane |

**Fig. 5 Multi-modal Captioning Results** with ImageBind-LLM. The training data is a collection of image-text pairs. We only train the parameters of the bind network during this stage.

ImageBind-LLM are able to connect other modalities with LLMs without any retraining. Although the pre-trained ImageBind-LLM cannot generate long and detailed description about the input, we believe that the basic ability acquired at this stage is very important for subsequent instruction tuning.

**Multi-modality Instruction Tuning.** By fine-tuning ImageBind-LLM on the mixture of language instruction data and visual instruction data, we equip ImageBind-LLM with both language and visual instruction-following abilities. More importantly, we can enter free-form instructions for modals other than images. Taking Fig. 6 as an example, ImageBind-LLM can not only answer questions about images, but also follow instructions of other modalities. We can ask the model to write a poem based on an audio (4th item in Fig. 6), introduce the content of a movie (5th item in Fig. 6) and even describe a 3D object (7-8th item in Fig. 6).

### 4.4.2 Mixture of Modalities

Since ImageBind-LLM unifies different modalities in one feature space, we can mix different modalities as the final input. In detail, we first employ modality encoders to extract modality tokens. We then mix the modality tokens by a set of coefficients. Finally, mixed tokens are fed into LLMs for multi-modality reasoning and understanding. For example, mixing images and audios allows the

model to better understand the current scene (1-2th item in Fig. 12); The combination of image and point cloud enrich the content of the picture (3-th item in Fig. 7); Using video frames and corresponding audio as input facilitates video understanding.

### 4.5 Advanced Applications

**Bilingual ImageBind-LLM.** By joint fine-tuning the model on both English and Chinese instruction data, we turn ImageBind-LLM into a bilingual model. As shown in Fig. 8, ImageBind-LLM can generate accurate Chinese response for a Chinese instruction. However, since our Chinese visual instruction capacity is emerged in a zero-shot manner, it usually fails to generate long and detailed Chinese responses. We believe that through a stronger bilingual LLM and collecting Chinese visual instruction data, its Chinese instruction capacity can be improved in the future.
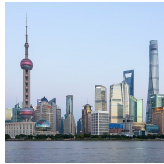
**Integration with Stable Diffusion.** With LLMs, we realize any-to-language decoding. Similarly, we employ Stable Diffusion for any-to image decoding. As shown in Fig. 8, ImageBind-LLM generates both language and image responses for a give input (*e.g.*, image, audio and point cloud). Compared with language response, the generated image carries more details about the input. For example, an image generated from a car horn can clearly depict the scene and its main elements.

13

**Image**

**Instruction:** Tell me about this game.
**Answer:** The game is called "The Legend of Zelda: Breath of the Wild," which is a popular action-adventure game developed by Nintendo for the Nintendo Switch and the Wii U console. The game is set in the fictional kingdom of Hyrule and follows the story of Link, a young boy who embarks on a quest to save the kingdom and rescue Princess Zelda. The game features open-world exploration, puzzle-solving, and various combat mechanics, making it a highly engaging and immersive experience for players.

**Instruction:** What city was this picture taken in and why?
**Answer:** The picture was taken in Shanghai, China. The presence of the tall buildings, the river, and the skyline suggest that the image was captured in a large city with a prominent urban landscape. Shanghai is known for its modern architecture, such as the Oriental Pearl Tower, and its proximity to the Huangpu River, which is a popular tourist attraction. The city's skyline, with its tall buildings and the river, creates a visually appealing and iconic backdrop for the image.

**Audio**

**Instruction:** How do you feel when you hear this music?
**Answer:** When I hear this music, I feel a sense of adventure and excitement. The combination of the electronic beats, the space-themed imagery, and the red and black color scheme create a thrilling atmosphere that invites listeners to explore new worlds and embark on a journey of discovery. The music may also evoke feelings of freedom and independence, as it may be associated with outer space or other far-off destinations.
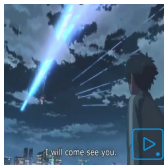
(electronic music)

**Instruction:** Write a poem based on this audio.
**Answer:** Amidst the crashing waves, a lone bird flies; Soaring high above the ocean's blue sky; The waves crashing, the winds blowing; A symphony of nature, a sight to see.

(sea wave)

**Video**

**Instruction:** Introduce me this movie.
**Answer:** The movie is titled "Your Name" and is a Japanese animated film released in 2016. It is a heartwarming and emotional story about two teenagers who discover a connection between their lives, and their journey to find the truth behind their shared experience. The film has received widespread critical acclaim and has been a commercial success, winning numerous awards and accolades, including the Academy Award for Best Animated Feature Film. The movie is known for its beautiful animation, captivating storytelling, and thought-provoking themes.

**Instruction:** Do you find this sport exciting? Describe this video.
**Answer:** Yes, this sport is exciting, as it showcases the thrilling experience of skydiving while also providing a unique perspective of the landscape below. The video captures the skydivers in mid-air, performing tricks and enjoying the adrenaline rush of the sport. The image of the parachutist flying through the sky with the landscape below is visually stunning and captures the essence of the freedom and adventure associated with skydiving.

**Point Cloud**

**Instruction:** Describe the 3D object in detail.
**Answer:** The 3D object is a grand piano, which is a musical instrument with a keyboard and a set of strings. The piano has a piano keyboard, which consists of white and black keys arranged in a grid-like pattern. The keys are connected to hammers that strike the strings, creating the sound of the piano. The piano in the image is a small, two-tiered model, with a brown and white color scheme. The piano is sitting on top of a bench, which is a standard accompaniment for a piano. The bench has a supportive design that allows the player to sit comfortably while playing the instrument.
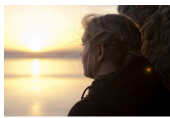
**Instruction:** Describe the 3D object in detail.
**Answer:** The 3D object is a flower pot with a plant inside it. The plant is a small, green, and leafy plant, which is growing out of the pot. The pot is made of clay and has a brown color. The plant is positioned towards the front of the pot, with its leaves spreading outwards. The pot is placed on a table, and the close-up view of the plant and pot showcases the details of the plant and the pot's design.
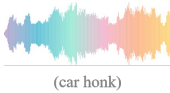
**Fig. 6 Multi-modal Instruction Examples with ImageBind-LLM.**

14

**Image+Audio**

**Instruction:** Guess the girl's mood based on the background music and explain the reason.
**Answer:** The girl in the image appears to be in a contemplative or melancholic mood, as the background music is a sad piano melody. This choice of music suggests that the scene is intended to evoke a sense of introspection, solitude, or longing. The combination of the girl's attire, the lone tree, and the melancholic music creates a mood of introspection and contemplation. The use of such music can help to enhance the emotional impact of the scene and provide a deeper understanding of the character's thoughts and feelings.

(sad piano)

**Instruction:** Describe the scene based on the input picture and audio.
**Answer:** The scene depicted in the image is a busy city street filled with traffic. There are numerous cars, trucks, and buses driving down the street, creating a bustling atmosphere. The street is crowded with vehicles, and the traffic appears to be moving slowly. The audio accompanying the image suggests that the traffic is quite congested, with honking and other noises typical of a busy city street. The image captures the essence of urban life, showcasing the hustle and bustle of daily city life.

(car honk)

**Image+Point Cloud**

**Instruction:** Describe this scene.
**Answer:** The image features a blue car driving down a tree-lined road, surrounded by a lush green forest. The car is in the middle of the road, and it appears to be a compact or small vehicle. There are several trees on both sides of the road, creating a serene and peaceful atmosphere. The car is moving along the road, possibly enjoying the scenic beauty of the forest.

**Video+Audio**

**Instruction:** What is the man doing in this video? What do people think of him based on the audio?
**Answer:** The man is standing at a podium, giving a speech or presentation. People in the audience are listening to him attentively, and they are clapping and showing their appreciation for his words.

(applause)

**Fig. 7 Mixture of Modalities with ImageBind-LLM**.

**Integration with Modality Experts.** ImageBind-LLM is a general multi-modality LLM, but it still lacks fine-grained domain-specific knowledge. Therefore, integrating with domain experts is a feasible way to improve ImageBind-LLM's multi-modality understanding capacity. In Fig. 9, we use the output of an speech-to-text expert as context to enhance ImageBind-LLM's audio understanding capacity. We also show that class labels extracted by a video object detector can help the model catch details of a video clip.

**ImageBind-LLM with Object Detection Cascade.** By cascading ImageBind-LLM with object detectors, we can associate text response with objects in the image, which helps improve the reliability of multi-modal LLMs. As shown in Fig. 10, we can evaluate the reliability of ImageBind-LLM through detection results. At the same time, ImageBind-LLM makes object detectors also has the ability to reasoning.

**ImageBind-LLM as Chatbot.** Leveraging multi-turn conversation data for training, ImageBind-LLM showcases impressive visual understanding abilities and multimodal conversation skills. As shown in Fig. 11, ImageBind-LLM excels in these capabilities. Furthermore, it exhibits the capacity to comprehend and reason with multimodal content in response to user queries, making it a viable candidate for a multi-modality chatbot.

**ImageBind-LLM for API control.** By incorporating tool-usage knowledge into LLMs, ImageBind-LLM can effectively learn to invoke API controls, enabling it to tackle various visual tasks. As illustrated in Fig. 12, the finetuned ImageBind-LLM demonstrates its proficiency in accomplishing diverse visual tasks, including but not limited to object detection and image editing, by effectively employing different APIs. These results highlight the potential of the proposed
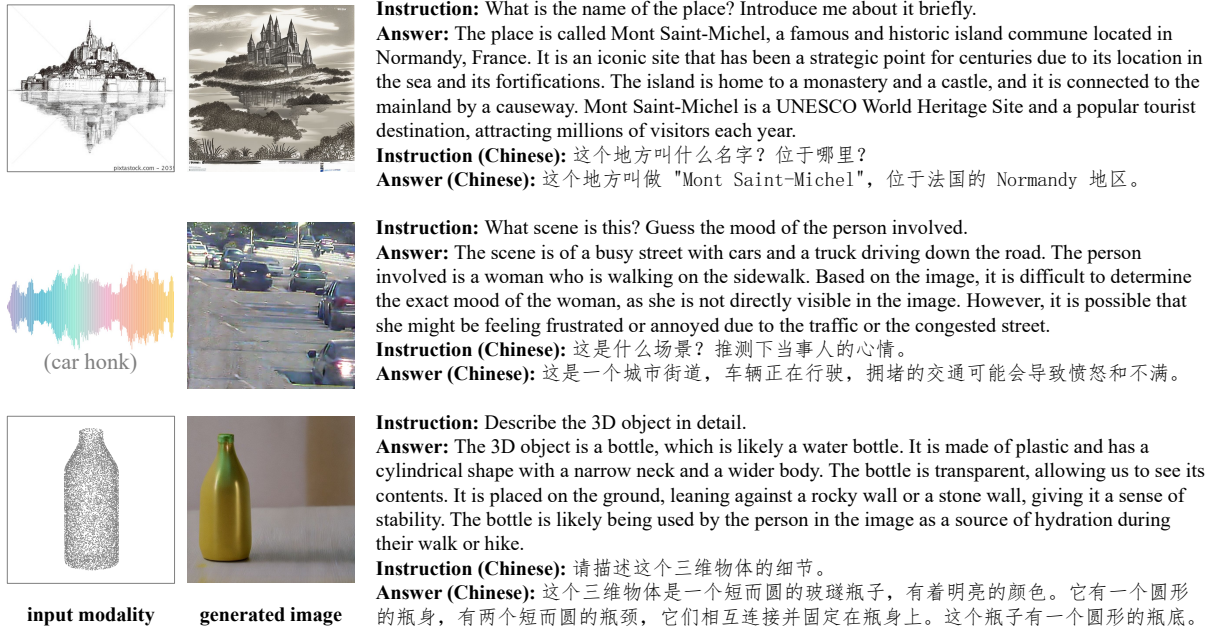
**Instruction:** What is the name of the place? Introduce me about it briefly.

**Answer:** The place is called Mont-Saint-Michel, a famous and historic island commune located in Normandy, France. It is an iconic site that has been a strategic point for centuries due to its location in the sea and its fortifications. The island is home to a monastery and a castle, and it is connected to the mainland by a causeway. Mont Saint-Michel is a UNESCO World Heritage Site and a popular tourist destination, attracting millions of visitors each year.

**Instruction (Chinese):** 这个地方叫什么名字？位于哪里？

**Answer (Chinese):** 这个地方叫做 "Mont Saint-Michel"，位于法国的 Normandy 地区。

**Instruction:** What scene is this? Guess the mood of the person involved.

**Answer:** The scene is of a busy street with cars and a truck driving down the road. The person involved is a woman who is walking on the sidewalk. Based on the image, it is difficult to determine the exact mood of the woman, as she is not directly visible in the image. However, it is possible that she might be feeling frustrated or annoyed due to the traffic or the congested street.

**Instruction (Chinese):** 这是什么场景？推测下当事人的心情。

**Answer (Chinese):** 这是一个城市街道，车辆正在行驶，拥堵的交通可能会导致愤怒和不满。

**Instruction:** Describe the 3D object in detail.

**Answer:** The 3D object is a bottle, which is likely a water bottle. It is made of plastic and has a cylindrical shape with a narrow neck and a wider body. The bottle is transparent, allowing us to see its contents. It is placed on the ground, leaning against a rocky wall or a stone wall, giving it a sense of stability. The bottle is likely being used by the person in the image as a source of hydration during their walk or hike.

**Instruction (Chinese):** 请描述这个三维物体的细节。

**Answer (Chinese):** 这个三维物体是一个短而圆的玻璃瓶子，有着明亮的颜色。它有一个圆形的瓶身，有两个短而圆的瓶颈，它们相互连接并固定在瓶身上。这个瓶子有一个圆形的瓶底。

**input modality**     **generated image**

**Fig. 8 Advanced Examples with Bilingual ImageBind-LLM and Integration with Stable Diffusion.**

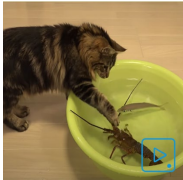*Speech-to-Text*

(Lover by Taylor Swift)

**Instruction:** What does this song convey?

**w/o expert:** The audio conveys a message that encourages viewers to "believe in the power of love."

**w/ Speech-to-text expert:** (**recognized text:** And ah, take me out, and take me home forever and ever. You're my my my my Lover.) The audio conveys a powerful and emotional message, possibly expressing a woman's desire to be with her lover forever. The lyrics suggest a strong, long-lasting bond between the two, and the woman's willingness to commit to this relationship. The song is likely to evoke feelings of love, passion, and devotion, as it is a romantic ballad.

*Video Object Detection*

**Instruction:** What is this cat doing?

**w/o expert:** The cat is standing on a table and watching the two fish swimming in a bowl.

**w/ Video Object Detection expert:** (**detected object classes:** cat, small fish, shrimp, water basin, floor) The cat is standing on the floor next to a water basin, and it appears to be observing or possibly playing with a small fish or shrimp that is swimming in the water.
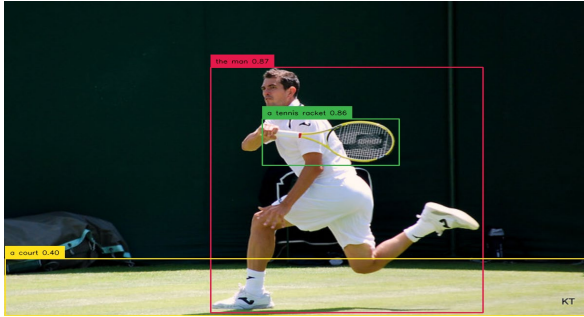
**Fig. 9 Advanced Examples by Integrating ImageBind-LLM with Modality Experts**.

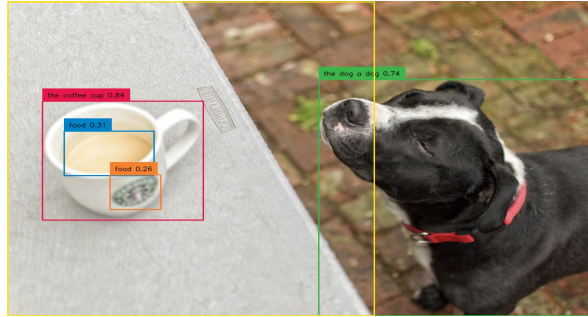ImageBind-LLM in the context of multi-modal tool usage.

## 4.6 Failure Cases

Although ImageBind-LLM can handle multi-modality inputs simultaneously, it is not perfect. As discussed in the previous quantitative analysis, ImageBind-LLM is weak compared to other VLMs. Firstly, ImageBind-LLM often suffers from hallucination issues for descriptive instructions. As shown in Fig. 13, ImageBind-LLM tends to describe objects not shown in the image. There maybe two possible reasons: (1) ImageBind-LLM only injects one global visual token into LLMs, which is much smaller than other models (10 for LLaMA-Adapter, 32 for MiniGPT4 and 256 for LLaVA). Therefore, LLM cannot get enough visual information in the self-Attention layers. (2)

**Instruction:** What is the man doing?
**Answer:** The man is playing tennis on a court, holding a tennis racket, and getting ready to hit the ball.
**Nouns:** The man, tennis, a court, a tennis racket, the ball.

**Instruction:** What is the dog want to do?
**Answer:** The dog wants to eat the coffee cup, likely because it is a dog and is attracted to food, especially when it is placed on a surface like a table.
**Nouns:** The dog, the coffee cup, food, a table.

**Fig. 10** Advanced Examples with Object Detection Cascade.

The quality of instruction tuning data is not high enough. For example, the visual instruction data from LLaVA is all generated by vision experts and GPT4, lacking human checks and corrections. Therefore, we will build a human-verified high-quality dataset in the future. Secondly, modality mixup is fragile when the two modalities represent different concepts. Fig. 13 gives an example of mixing a "sky" image and a "dog barking" audio. We can see that ImageBind-LLM is sensitive to the modality mixing ratio.

## 5 Conclusion

In this work, we propose to tune LLMs into multi-modality instruction models with the help of ImageBind, named ImageBind-LLM. In contrast to prior language instruction models and image instruction models, ImageBind-LLM unifies image, audio, 3D point clouds and video into one multi-modality LLM. We achieves this by simply aligning ImageBind's visual encoder with an LLM via a learnable bind network. Thanks to the binding property of ImageBind, we can directly feed multi-modality inputs into ImageBind-LLM for inference without any training. We also propose a training-free image cache model to mitigate the modality discrepancy between training and inference. We evaluate ImageBind-LLM across 27 traditional vision-language datasets and a new multimodal LLM benchmark MME, where ImageBind-LLM achieves comparable performance with recent mutlimodal LLMs. We also

give extensive qualitative analysis to demonstrate ImageInd-LLM's multi-modality understanding capacity. In the future, we plan to enhance ImageBind-LLM by increasing the number of multi-modality tokens. Besides, integrating more modalities into ImageBind-LLM is also a promising approach to enhance its multimodal capacity.

## References

[1] Gao, P., Han, J., Zhang, R., Lin, Z., Geng, S., Zhou, A., Zhang, W., Lu, P., He, C., Yue, X., Li, H., Qiao, Y.: LLaMA-Adapter V2: Parameter-Efficient Visual Instruction Model (2023)

[2] Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. arXiv preprint arXiv:2304.08485 (2023)

[3] Zhu, D., Chen, J., Shen, X., Li, X., Elhoseiny, M.: Minigpt-4: Enhancing vision-language understanding with advanced large language models. arXiv preprint arXiv:2304.10592 (2023)

[4] ChatGPT. https://chat.openai.com (2023)

[5] OpenAI: Gpt-4 technical report. ArXiv **abs/2303.08774** (2023)

[6] Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., Hashimoto,

**Image**

Introduce me this painting.

The painting is a famous work of art titled "Starry Night" by Vincent van Gogh. It was created in 1889 and is one of the most recognizable and iconic paintings in the world. The painting depicts a night sky with a swirling mass of stars, a crescent moon, and a small village in the distance. The artist's use of bold, vibrant colors and his unique brushstrokes have made this painting a timeless masterpiece. The painting is a testament to van Gogh's passion for art and his ability to capture the beauty and wonder of the night sky.

Tell me more about his work.

Vincent van Gogh was a Dutch post-impressionist painter who is widely recognized for his unique style and the emotional depth of his paintings. He is known for his masterpiece "Starry Night," but he also created many other iconic works, including "The Café Terrace at Night," "The Wheat Field with Starry Sky," "The Bedroom," "The Sower," "The Potato Eaters," "The Starry Night Over the Rhone," and "Irises," among others.
Van Gogh's artistic career spanned only ten years, from 1880 to 1890, during which time he produced over 2,000 paintings and drawings. He is best known for his series of paintings, including "Starry Night," "Irises," "The Wheat Field with the Starry Sky," and "The Café Terrace at Night."
Van Gogh's artistic style has been greatly admired and influential, and his work has been the subject of numerous exhibitions and studies. He is considered one of the greatest artists of all time, and his paintings continue to inspire and captivate art lovers and collectors alike.

**Video**

**Audio**



Describe the content in the video.

The video shows the two young girls sitting on the grass, likely taking a break from walking or enjoying their time together. They are surrounded by the same grassy field or meadow as before. The scene appears to be a peaceful and relaxing moment shared between the two girls.

What is the weather like?

The weather in the image appears to be sunny and pleasant, as the girls are sitting outside in the sunlight.

Guess the mood of the girls based on the audio.

Based on the audio, the mood of the girls appears to be joyful and relaxed. They are likely having a good time together, enjoying each other's company and the beautiful outdoor setting.

**Fig. 11 Advanced Examples with ImageBind-LLM as Chatbot.**

T.B.: Stanford Alpaca: An Instruction-following LLaMA model. GitHub (2023)

[7] Zhang, R., Han, J., Zhou, A., Hu, X., Yan, S., Lu, P., Li, H., Gao, P., Qiao, Y.: Llama-adapter: Efficient fine-tuning of language models with zero-init attention. arXiv preprint arXiv:2303.16199 (2023)
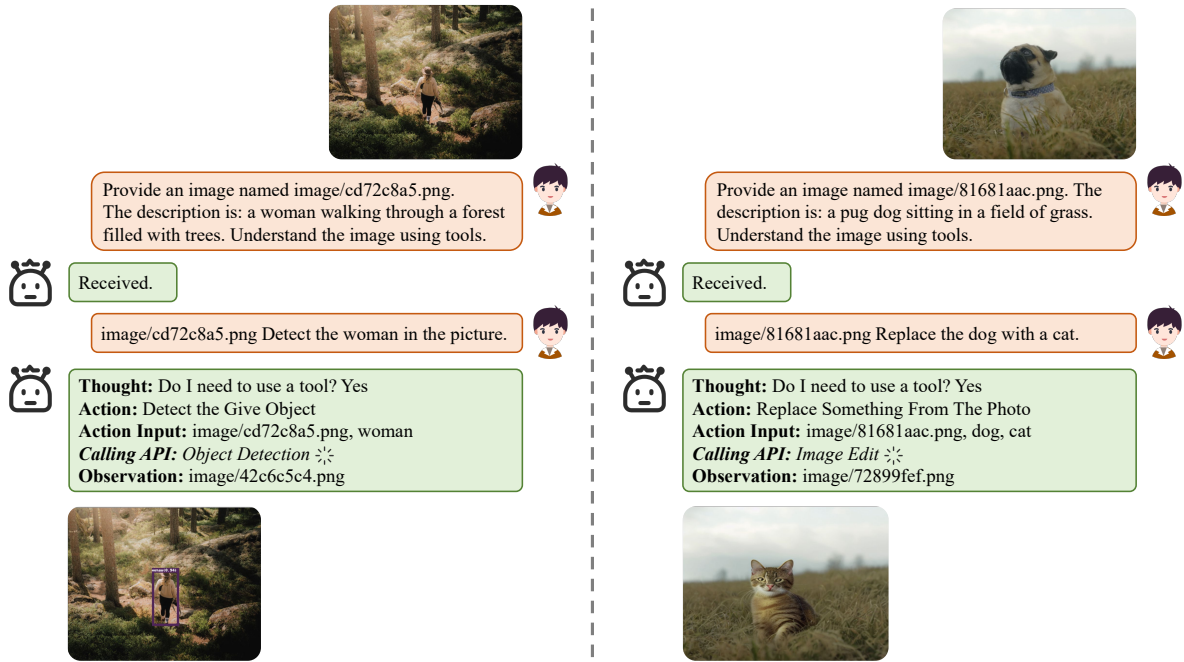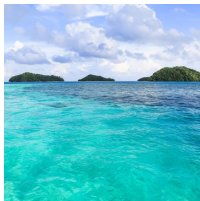
[8] Alpaca-LoRA. https://github.com/tloen/alpaca-lora (2023)

**Fig. 12  Advanced Examples with ImageBind-LLM for API Control.**

*Hallucination*



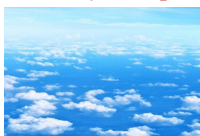**Instruction:** Describe this image in detail.
**Answer:** The image features a lighthouse standing tall on a rocky hill, overlooking the ocean. The lighthouse is white and has a red top, making it a prominent landmark. The sky in the background is cloudy, adding a dramatic effect to the scene. There are several people scattered around the area, likely enjoying the view of the lighthouse and the ocean. Some of them are closer to the lighthouse, while others are further away, creating a sense of depth in the scene.



**Instruction:** Describe this image in detail.
**Answer:** The image features a large body of water with a beautiful blue color, where several small boats are floating on the surface. There are at least 11 boats visible in the scene, scattered across the water. Some of the boats are closer to the foreground, while others are further away, creating a sense of depth in the scene. The boats are of various sizes and shapes, adding to the diversity of the scene. The water appears calm and serene, making it an ideal setting for a day of boating and relaxation.

*Modality Mixup is Fragile*



(dog barking)

**Instruction:** Describe this scene briefly.
**Image*0.5+Audio*0.5: (Missing audio keyword "dog")** The scene features a blue sky with white, fluffy clouds, creating a beautiful backdrop for the image. The sky is filled with numerous small, white clouds that are scattered throughout the sky, giving it a serene and peaceful appearance.
**Image*0.4+Audio*0.6: (Include both "sky" and "dog")** The scene features a dog standing on a sidewalk, looking up at the sky. The dog appears to be looking at a person's shadow, which is cast on the ground. The dog is wearing a collar, and the sidewalk is located near a fence.
**Image*0.2+Audio*0.8: (Missing image keyword "sky")** The scene features a dog standing on a sidewalk, looking at a person who is walking a cat on a leash. The dog appears to be curious about the cat, as it is staring at the cat and the person walking it.
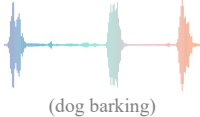
**Fig. 13  Failure Cases.**

19

[9] Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J.E., Stoica, I., Xing, E.P.: Vicuna: An Open-Source Chatbot Impressing GPT-4 with 90%* ChatGPT Quality. https://lmsys.org/blog/2023-03-30-vicuna/ (2023)

[10] Peng, B., Li, C., He, P., Galley, M., Gao, J.: Instruction tuning with gpt-4. arXiv preprint arXiv:2304.03277 (2023)

[11] Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., Lample, G.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023)

[12] Girdhar, R., El-Nouby, A., Liu, Z., Singh, M., Alwala, K.V., Joulin, A., Misra, I.: Imagebind: One embedding space to bind them all. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 15180–15190 (2023)

[13] Guo, Z., Zhang, R., Zhu, X., Tang, Y., Ma, X., Han, J., Chen, K., Gao, P., Li, X., Li, H., et al.: Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. arXiv preprint arXiv:2309.00615 (2023)

[14] Zhang, R., Fang, R., Gao, P., Zhang, W., Li, K., Dai, J., Qiao, Y., Li, H.: Tip-adapter: Training-free clip-adapter for better vision-language modeling. arXiv preprint arXiv:2111.03930 (2021)

[15] Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021)

[16] Xie, E., Yao, L., Shi, H., Liu, Z., Zhou, D., Liu, Z., Li, J., Li, Z.: Difffit: Unlocking transferability of large diffusion models via simple parameter-efficient fine-tuning. arXiv preprint arXiv:2304.06648 (2023)

[17] Zaken, E.B., Ravfogel, S., Goldberg, Y.: Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. arXiv preprint arXiv:2106.10199 (2021)

[18] Frankle, J., Schwab, D.J., Morcos, A.S.: Training batchnorm and only batchnorm: On the expressive power of random features in cnns. arXiv preprint arXiv:2003.00152 (2020)

[19] Giannou, A., Rajput, S., Papailiopoulos, D.: The expressive power of tuning only the norm layers. arXiv preprint arXiv:2302.07937 (2023)

[20] Lu, P., Mishra, S., Xia, T., Qiu, L., Chang, K.-W., Zhu, S.-C., Tafjord, O., Clark, P., Kalyan, A.: Learn to explain: Multimodal reasoning via thought chains for science question answering. In: The 36th Conference on Neural Information Processing Systems (NeurIPS) (2022)

[21] OpenAI: GPT-4 Technical Report (2023)

[22] Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual Instruction Tuning. arXiv:2304.08485 (2023)

[23] Li, B., Zhang, Y., Chen, L., Wang, J., Yang, J., Liu, Z.: Otter: A multi-modal model with in-context instruction tuning. arXiv preprint arXiv:2305.03726 (2023)

[24] Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597 (2023)

[25] Chen, G., Zheng, Y.-D., Wang, J., Xu, J., Huang, Y., Pan, J., Wang, Y., Wang, Y., Qiao, Y., Lu, T., et al.: Videollm: Modeling video sequence with large language models. arXiv preprint arXiv:2305.13292 (2023)

[26] Zhang, H., Li, X., Bing, L.: Video-llama: An instruction-tuned audio-visual language model for video understanding. arXiv preprint arXiv:2306.02858 (2023)

[27] Radford, A., Kim, J.W., Hallacy, C., Ramesh,

A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., *et al.*: Learning transferable visual models from natural language supervision. In: International Conference on Machine Learning, pp. 8748–8763 (2021). PMLR

[28] Jia, C., Yang, Y., Xia, Y., Chen, Y.-T., Parekh, Z., Pham, H., Le, Q., Sung, Y.-H., Li, Z., Duerig, T.: Scaling up visual and vision-language representation learning with noisy text supervision. In: International Conference on Machine Learning, pp. 4904–4916 (2021). PMLR

[29] Yuan, L., Chen, D., Chen, Y.-L., Codella, N., Dai, X., Gao, J., Hu, H., Huang, X., Li, B., Li, C., et al.: Florence: A new foundation model for computer vision. arXiv preprint arXiv:2111.11432 (2021)

[30] Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., *et al.*: Flamingo: a visual language model for few-shot learning. Advances in Neural Information Processing Systems **35**, 23716–23736 (2022)

[31] Zhang, W., Shi, H., Guo, J., Zhang, S., Cai, Q., Li, J., Luo, S., Zhuang, Y.: Magic: Multimodal relational graph adversarial inference for diverse and unpaired text-based image captioning. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 36, pp. 3335–3343 (2022)

[32] Guzhov, A., Raue, F., Hees, J., Dengel, A.: Audioclip: Extending clip to image, text and audio. In: ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 976–980 (2022). IEEE

[33] Zhang, R., Guo, Z., Zhang, W., Li, K., Miao, X., Cui, B., Qiao, Y., Gao, P., Li, H.: Pointclip: Point cloud understanding by clip. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8552–8562 (2022)

[34] Su, Y., Lan, T., Li, H., Xu, J., Wang, Y., Cai, D.: Pandagpt: One model to instruction-follow them all. arXiv preprint arXiv:2305.16355 (2023)

[35] Khandelwal, U., Levy, O., Jurafsky, D., Zettlemoyer, L., Lewis, M.: Generalization through memorization: Nearest neighbor language models. arXiv preprint arXiv:1911.00172 (2019)

[36] Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al.: Matching networks for one shot learning. Advances in neural information processing systems **29** (2016)

[37] Zhang, R., Wang, L., Wang, Y., Gao, P., Li, H., Shi, J.: Parameter is not all you need: Starting from non-parametric networks for 3d point cloud analysis. CVPR 2023 (2023)

[38] Udandarao, V., Gupta, A., Albanie, S.: Sus-x: Training-free name-only transfer of vision-language models. arXiv preprint arXiv:2211.16198 (2022)

[39] Zhu, X., Zhang, R., He, B., Zhou, A., Wang, D., Zhao, B., Gao, P.: Not all features matter: Enhancing few-shot clip with adaptive prior refinement. arXiv preprint arXiv:2304.01195 (2023)

[40] Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., *et al.*: Laion-5b: An open large-scale dataset for training next generation image-text models. Advances in Neural Information Processing Systems **35**, 25278–25294 (2022)

[41] Sharma, P., Ding, N., Goodman, S., Soricut, R.: Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 2556–2565 (2018)

[42] Changpinyo, S., Sharma, P., Ding, N., Soricut, R.: Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3558–3568

(2021)

[43] Zhang, B., Sennrich, R.: Root mean square layer normalization. Advances in Neural Information Processing Systems **32** (2019)

[44] Hendrycks, D., Gimpel, K.: Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415 (2016)

[45] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778 (2016)

[46] https://sharegpt.com/

[47] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10684–10695 (2022)

[48] Huang, R., Huang, J., Yang, D., Ren, Y., Liu, L., Li, M., Ye, Z., Liu, J., Yin, X., Zhao, Z.: Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models. arXiv preprint arXiv:2301.12661 (2023)

[49] Sanghi, A., Chu, H., Lambourne, J.G., Wang, Y., Cheng, C.-Y., Fumero, M., Malekshan, K.R.: Clip-forge: Towards zero-shot text-to-shape generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 18603–18613 (2022)

[50] Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., Li, C., Yang, J., Su, H., Zhu, J., et al.: Grounding dino: Marrying dino with grounded pre-training for open-set object detection. arXiv preprint arXiv:2303.05499 (2023)

[51] Bird, S., Klein, E., Loper, E.: Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit. ” O'Reilly Media, Inc.”, ??? (2009)

[52] Yang, R., Song, L., Li, Y., Zhao, S., Ge, Y., Li, X., Shan, Y.: Gpt4tools: Teaching large language model to use tools via self-instruction. arXiv preprint arXiv:2305.18752 (2023)

[53] Chen, X., Fang, H., Lin, T.-Y., Vedantam, R., Gupta, S., Dollár, P., Zitnick, C.L.: Microsoft coco captions: Data collection and evaluation server. arXiv preprint arXiv:1504.00325 (2015)

[54] Ordonez, V., Kulkarni, G., Berg, T.L.: Im2text: Describing images using 1 million captioned photographs. In: Neural Information Processing Systems (NIPS) (2011)

[55] Byeon, M., Park, B., Kim, H., Lee, S., Baek, W., Kim, S.: COYO-700M: Image-Text Pair Dataset. https://github.com/kakaobrain/coyo-dataset (2022)

[56] Zhu, W., Hessel, J., Awadalla, A., Gadre, S.Y., Dodge, J., Fang, A., Yu, Y., Schmidt, L., Wang, W.Y., Choi, Y.: Multimodal C4: An open, billion-scale corpus of images interleaved with text. arXiv preprint arXiv:2304.06939 (2023)

[57] Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: ICML (2022)

[58] Johnson, J., Douze, M., Jégou, H.: Billion-scale similarity search with GPUs. IEEE Transactions on Big Data **7**(3), 535–547 (2019)

[59] Fu, C., Chen, P., Shen, Y., Qin, Y., Zhang, M., Lin, X., Qiu, Z., Lin, W., Yang, J., Zheng, X., et al.: Mme: A comprehensive evaluation benchmark for multimodal large language models. arXiv preprint arXiv:2306.13394 (2023)

[60] Dai, W., Li, J., Li, D., Tiong, A.M.H., Zhao, J., Wang, W., Li, B., Fung, P., Hoi, S.: Instructblip: Towards general-purpose vision-language models with instruction tuning. arXiv preprint arXiv:2305.06500 (2023)

[61] Mishra, A., Alahari, K., Jawahar, C.V.: Top-down and bottom-up cues for scene text recognition. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition, pp. 2687–2694 (2012). https://doi.org/10.1109/CVPR.2012.6247990

[62] Karatzas, D., Shafait, F., Uchida, S., Iwamura, M., Bigorda, L.G.i., Mestre, S.R., Mas, J., Mota, D.F., Almazàn, J.A., Heras, L.P.: Icdar 2013 robust reading competition. In: 2013 12th International Conference on Document Analysis and Recognition, pp. 1484–1493 (2013). https://doi.org/10.1109/ICDAR.2013.221

[63] Karatzas, D., Gomez-Bigorda, L., Nicolaou, A., Ghosh, S., Bagdanov, A., Iwamura, M., Matas, J., Neumann, L., Chandrasekhar, V.R., Lu, S., Shafait, F., Uchida, S., Valveny, E.: Icdar 2015 competition on robust reading. In: 2015 13th International Conference on Document Analysis and Recognition (ICDAR), pp. 1156–1160 (2015). https://doi.org/10.1109/ICDAR.2015.7333942

[64] Ch'ng, C.K., Chan, C.S.: Total-text: A comprehensive dataset for scene text detection and recognition. In: 2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR), vol. 01, pp. 935–942 (2017). https://doi.org/10.1109/ICDAR.2017.157

[65] Risnumawan, A., Shivakumara, P., Chan, C.S., Tan, C.L.: A robust arbitrary text detection system for natural scene images. Expert Systems with Applications **41**(18), 8027–8048 (2014) https://doi.org/10.1016/j.eswa.2014.07.008

[66] Shi, C., Wang, C., Xiao, B., Gao, S., Hu, J.: End-to-end scene text recognition using tree-structured models. Pattern Recognition **47**(9), 2853–2866 (2014) https://doi.org/10.1016/j.patcog.2014.03.023

[67] Phan, T.Q., Shivakumara, P., Tian, S., Tan, C.L.: Recognizing text with perspective distortion in natural scenes. In: 2013 IEEE International Conference on Computer Vision, pp. 569–576 (2013). https://doi.org/10.1109/ICCV.2013.76

[68] Veit, A., Matera, T., Neumann, L., Matas, J., Belongie, S.J.: Coco-text: Dataset and benchmark for text detection and recognition in natural images. ArXiv **abs/1601.07140** (2016)

[69] Xie, X., Fu, L., Zhang, Z., Wang, Z., Bai, X.: Toward understanding wordart: Corner-guided transformer for scene text recognition (2022)

[70] Liu, Y., Jin, L., Zhang, S., Luo, C., Zhang, S.: Curved scene text detection via transverse and longitudinal sequence connection. Pattern Recogn. **90**(C), 337–345 (2019) https://doi.org/10.1016/j.patcog.2019.02.002

[71] Wang, Y., Xie, H., Fang, S., Wang, J., Zhu, S., Zhang, Y.: From two to one: A new scene text recognizer with visual language modeling network. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 14194–14203 (2021)

[72] Huang, Z., Chen, K., He, J., Bai, X., Karatzas, D., Lu, S., Jawahar, C.V.: Icdar2019 competition on scanned receipt ocr and information extraction. In: 2019 International Conference on Document Analysis and Recognition (ICDAR), pp. 1516–1520 (2019). https://doi.org/10.1109/ICDAR.2019.00244

[73] Mathew, M., Karatzas, D., Jawahar, C.: Docvqa: A dataset for vqa on document images. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 2200–2209 (2021)

[74] Singh, A., Natarajan, V., Shah, M., Jiang, Y., Chen, X., Batra, D., Parikh, D., Rohrbach, M.: Towards vqa models that can read. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 8309–8318 (2019). https://doi.org/10.1109/CVPR.2019.00851

[75] Furkan Biten, A., Tito, R., Mafla, A., Gomez, L., Rusiñol, M., Mathew, M., Jawahar, C.V., Valveny, E., Karatzas, D.: Icdar

23

2019 competition on scene text visual question answering. In: 2019 International Conference on Document Analysis and Recognition (ICDAR), pp. 1563–1570 (2019). https://doi.org/10.1109/ICDAR.2019.00251

[76] Mishra, A., Shekhar, S., Singh, A.K., Chakraborty, A.: Ocr-vqa: Visual question answering by reading text in images. In: 2019 International Conference on Document Analysis and Recognition (ICDAR), pp. 947–952 (2019). https://doi.org/10.1109/ICDAR.2019.00156

[77] Marino, K., Rastegari, M., Farhadi, A., Mottaghi, R.: Ok-vqa: A visual question answering benchmark requiring external knowledge. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 3190–3199 (2019). https://doi.org/10.1109/CVPR.2019.00331

[78] Hudson, D.A., Manning, C.D.: Gqa: A new dataset for real-world visual reasoning and compositional question answering. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 6693–6702 (2019). https://doi.org/10.1109/CVPR.2019.00686

[79] Lu, P., Qiu, L., Chen, J., Xia, T., Zhao, Y., Zhang, W., Yu, Z., Liang, X., Zhu, S.-C.: Iconqa: A new benchmark for abstract diagram understanding and visual language reasoning. In: The 35th Conference on Neural Information Processing Systems (NeurIPS) Track on Datasets and Benchmarks (2021)

[80] Liu, F., Emerson, G.E.T., Collier, N.: Visual spatial reasoning. Transactions of the Association for Computational Linguistics (2023)

[81] Das, A., Kottur, S., Gupta, K., Singh, A., Yadav, D., Moura, J.M.F., Parikh, D., Batra, D.: Visual Dialog. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)

[82] Bigham, J.P., Jayant, C., Ji, H., Little, G., Miller, A., Miller, R.C., Miller, R., Tatarowicz, A., White, B., White, S., *et al.*: Vizwiz: nearly real-time answers to visual questions. In: Proceedings of the 23nd Annual ACM Symposium on User Interface Software and Technology, pp. 333–342 (2010)

[83] Awadalla, A., Gao, I., Gardner, J., Hessel, J., Hanafy, Y., Zhu, W., Marathe, K., Bitton, Y., Gadre, S., Jitsev, J., Kornblith, S., Koh, P.W., Ilharco, G., Wortsman, M., Schmidt, L.: OpenFlamingo. Zenodo (2023). https://doi.org/10.5281/zenodo.7733589 . https://doi.org/10.5281/zenodo.7733589