

# Advances in Large Multi-Modal Models from the Perspective of Representation Space Extension: A Survey

Anonymous ACL submission

## Abstract

The success of large language models (LLMs) has attracted much focus on extending these models to multi-modal domains, giving rise to large multi-modal models (LMMs). Unlike existing reviews that focus on specific model frameworks or scenarios, this paper aims to provide an encyclopedic survey on LMMs from a general perspective, i.e. **representation space extension**. By systematically analyzing the input-output representations of existing LMMs, this paper summarizes the design of model architectures to align the constructed multi-modal representation space. Lastly, this paper demonstrates the extensibility of LMMs as embodied agents in view of proposed representation space extension. With the insights revealed through surveying the field, this paper discusses several fundamental problems of constructing LMMs and inspires future work at the end.

## 1 Introduction

The goal of AI research is to build versatile intelligent systems capable of fulfilling tasks across diverse scenarios. Recently, the generalization and interactivity demonstrated by large language models (LLMs) have significantly advanced the progress towards general-purpose AI (OpenAI, 2023; Touvron et al., 2023b; Bai et al., 2023a; AI@Meta, 2024). To adapt these capabilities to multi-modal contexts, research on large multi-modal models (LMMs) is emerging, aiming to extend the input and output representation space of the language-based interface to more modalities. As shown in Figure 1, to extend the input space, existing methods introduce discretely or continuously encoded modality representations into the text input and learn cross-modal alignment from multi-modal intertwined data, enabling LMMs to understand multi-modal information (Li et al., 2023b; Liu et al., 2024b; Bai et al., 2023b; Ma et al., 2024). Similarly, the output space can be divided into multiple subspaces

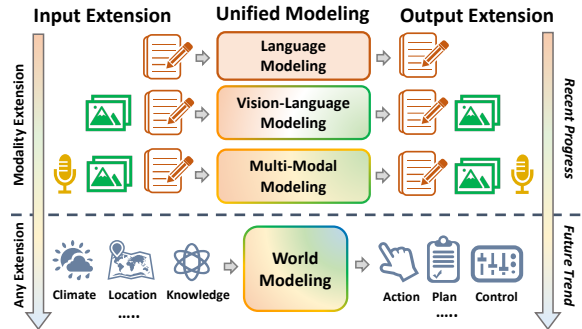


Figure 1: Illustration of the general LMM framework: expanding input and output representation space to more modalities and aligning representations across modalities through unified multi-modal modeling.

of different modalities, which are further aligned with corresponding modality decoders to generate multi-modal content (Koh et al., 2024; Zhang et al., 2024a; Wu et al., 2023d; Zhan et al., 2024a).

Although there are several surveys that detail the current progress in constructing LMMs, most of these works are limited to specific sub-problems in the construction of LMMs, such as applications in specific modalities (Tang et al., 2023b; Latif et al., 2023) and scenarios (Xiao et al., 2024; Cui et al., 2024). Meanwhile, most existing reviews focus on a specific type of model framework: encoding information from other modalities in a continuous manner and aligning them with text embeddings through connection modules (Wu et al., 2023c; Caffagni et al., 2024), neglecting related research on other architectures, such as unified discretely represented LMMs (Team, 2024; Zhan et al., 2024a). These limitations prevent existing reviews from adequately covering research problems in LMM construction and limit their applicability.

To this end, this survey aims to summarize related works from a more general perspective: **the extension of input-output representation space**. As illustrated in Figure 1, existing LMMs can be systematically summarized from this view, encompassing various modalities, scenarios, and model

068	architectures, while also leaving room for further	<b>2.1 Encode Input Representation</b>	116
069	exploration to more modalities and scenarios. We	Regarding the input, the core research problems	117
070	omit the details about data and evaluations that have	involve how to code the representations of each	118
071	been sufficiently reviewed by previous works (Li	modality and how to integrate them into a multi-	119
072	and Lu, 2024; Huang and Zhang, 2024; Bai et al.,	modal input space (see lower part in Figure 2).	120
073	2024), but keep a tight focus on the architectures.		
074	To conduct a holistic survey, we follow a top-	<b>2.1.1 Textual Representation</b>	121
075	down logic to break down the construction of	As a discrete signal, text is a sequence composed of	122
076	LMMs into several sub-problems, providing de-	characters. Following the practice of LLMs, LMMs	123
077	tailed discussions to offer insights to readers. Par-	typically utilize tokenizers, such as BPE (Sennrich	124
078	ticularly, we try to answer the following questions.	et al., 2015; Radford et al., 2019), WordPiece (Wu,	125
079	(i) How can modality signals be encoded using	2016), and Unigram (Kudo, 2018), to merge char-	126
080	discrete or continuous representations, and how to	acters into sub-word tokens. Ultimately, texts are	127
081	construct multi-modal representation spaces? (§2)	represented as sequences of discrete tokens.	128
082	(ii) How to design model architectures to align the		
083	constructed multi-modal representation space? (§3)	<b>2.1.2 Visual Representation</b>	129
084	(iii) How to extend the representation space to real	For visual signals with spatial-temporal informa-	130
085	scenarios, i.e. embodied agents? (§4) This fur-	tion, LMMs mainly employ pre-trained visual en-	131
086	ther demonstrates the extensibility of LMMs from	coders for representing images (videos) into con-	132
087	the perspective discussed in this paper. Finally, in	tinuous features or discrete codes. Figure 5 in Ap-	133
088	§5, we summarize the discussion on the questions	pendix A shows the evolution of visual encoders.	134
089	raised above, providing readers with key take-home	Commonly adopted architectures of visual	135
090	messages and an outlook on future research.	encoders can be divided into two categories:	136
091	In summary, our contributions are threefold:	convolution-based (He et al., 2016; Liu et al.,	137
		2022b) and vision-Transformer-based mod-	138
092	• Going beyond specific scenarios and model	els (Dosovitskiy et al., 2020; Liu et al., 2021).	139
093	framework, we review the current LMMs	Both methods encode images into continuous	140
094	from a general perspective of input-output rep-	2D feature maps. These continuous features can	141
095	resentation space extension.	be further compressed into discrete visual codes	142
096		through vector quantization (VQ) by learning a	143
097	• Based on the structure of input-output spaces,	fixed-size visual codebook (Van Den Oord et al.,	144
098	we systematically review the existing mod-	2017; Esser et al., 2021). In addition, models like	145
099	els, including mainstream models based on	Fuyu (Bavishi et al., 2023) do not rely on visual	146
100	discrete-continuous hybrid spaces and models	encoders and directly use pixel values of image	147
101	with unified multi-modal discrete representa-	patches as the visual representations.	148
102	tions. Furthermore, we summarize the design	Based on the sequence modeling framework	149
103	of model architectures to align the constructed	of current LMMs, multiple images can be intu-	150
104	multi-modal representation space.	itively arranged in the input sequence (Luo et al.,	151
105		2023b; Zhang et al., 2023b; Li et al., 2023a; Yu	152
106	• We elaborate on how to extend LMMs to em-	et al., 2024b). For videos, where images (frames)	153
107	body scenarios to highlight the extensibility	are temporally related, spatial-temporal encoders	154
108	of LMMs from the input-output extension per-	such as TimeSformer (Bertasius et al., 2021) and	155
	spective. To the best of our knowledge, this is	VideoSwin (Liu et al., 2022a) can be further used	156
	the first survey to include embodied LMMs.	for encoding (Li et al., 2023c; Xu et al., 2023).	157
109			
	<b>2 Representation Space Extension</b>	<b>2.1.3 Multi-Modal Representation</b>	158
110	In this section, we introduce prevalent solutions	As illustrated in the lower part Figure 2, there exist	159
111	to construct multi-modal representation space. As	two mainstream types of multi-modal input space.	160
112	illustrated in Figure 2, existing methods can be		
113	categorized based on different input-output space	<b>Type A: Hybrid Input Space</b> Text are repre-	161
114	structures, and the extension to other modalities	sented in a discrete form, while visual signals are	162
115	can be summarized in a similar manner.	encoded in continuous representations, preseving	163
		the complete visual information. However, due to	164

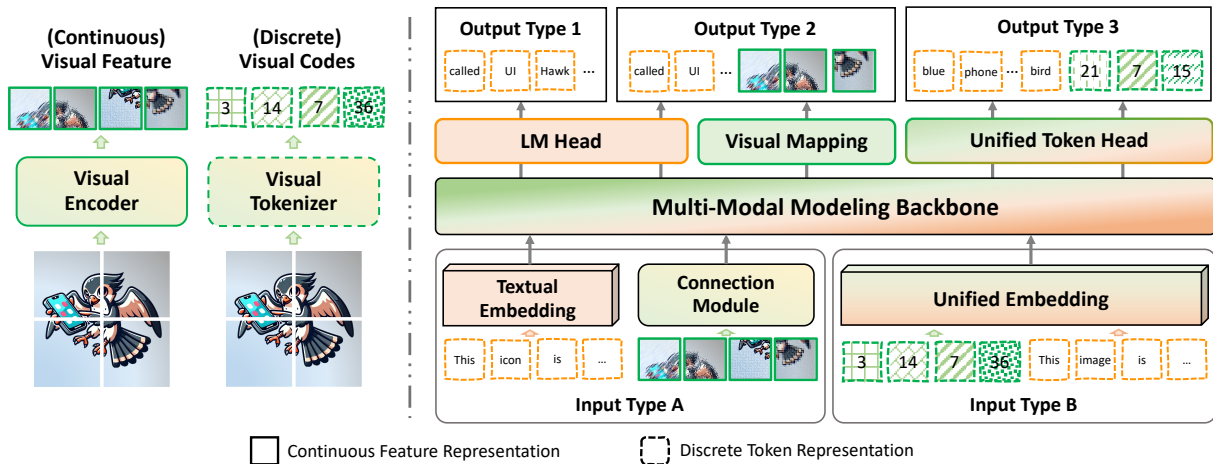


Figure 2: Summary and illustration of different input-output space structures for extension to vision modality.

the gap in the input space, connection modules are required to perform input-level cross-modal alignment, which is discussed in Section 3.

**Type B: Unified Discrete Input Space** Different from Type A, further quantizing visual representations into discrete visual codes facilitates the construction of a unified input space. A multi-model vocabulary can be intuitively integrated and directly used to support subsequent modeling.

### 2.1.4 Extension to More Input Modalities

Beyond the vision modality, signals from other modalities can be encoded and introduced into the input space following a similar paradigm. For example, various encoders can help encode audio into continuous (Hsu et al., 2021; Elizalde et al., 2023; Girdhar et al., 2023) or discrete (Zhang et al., 2023e) representations. As a step further, an arbitrary-modality input space can be represented in either hybrid (Wu et al., 2023d; Han et al., 2023; Tang et al., 2023c; Lu et al., 2023a) or unified discrete forms (Zhan et al., 2024a).

## 2.2 Decode Output Representation

Based on the input, backbones of LMMs present continuous multi-modal output representations which can be used to decode the output signals of different modalities. For example, with the commonly used causal modeling framework, the output representation can be leveraged to predict the signal at the next position in the sequence. Predicted token sequence can be converted to text with the tokenizer while different image generator can be adopted to decode images from output in different forms. In this section, we discuss the commonly adopted paradigms to partition the output space

of different modalities and perform corresponding decoding, as shown in the upper part of Figure 2.

### 2.2.1 Type 1: Text-Only Output Space

If only text output is required, similar to LLMs, discrete tokens can be generated from the output representations through a classification-based language modeling (LM) head and specific decoding strategies (Li et al., 2023b; Liu et al., 2024b). Please note that models that first generate text descriptions and then use external tools like Stable Diffusion and CLIP to generate or retrieve content in other modalities, such as Visual ChatGPT (Wu et al., 2023a), InternLM-XComposer series (Zhang et al., 2023c; Dong et al., 2024b), and Mini-Gemini (Li et al., 2024d), are also classified as text-only output models because they are not in an end-to-end manner.

### 2.2.2 Type 2: Hybrid Output Space

The hybrid output space includes the discrete text tokens and continuous visual features. Such output space is initially proposed to support image generation. A series of methods first introduce special tokens, such as the start and end tokens for images, or a series of consecutive placeholder tokens to indicate where images should be generated. The continuous visual representations at the corresponding positions are then connected to visual decoders (mainly Diffusion models (Rombach et al., 2022)) through visual mapping modules (Koh et al., 2024; Dong et al., 2024a; Zheng et al., 2024b; Sun et al., 2024b). Similar to the hybrid input space, visual mapping modules perform output-level alignment and typically requires further training.

### 2.2.3 Type 3: Unified Discrete Output Space

The unified discrete output space contains discrete text tokens and discrete visual codes. Based on the



joint vocabulary constructed within Type B input space described in Section 2.1.3, image generation is naturally integrated into the token decoding process. The predicted visual codes are fed to the corresponding codebook detokenizer to generate the image (Ge et al., 2023b; Team, 2024).

#### 2.2.4 Extension to More Output Modalities

Type 2 and Type 3 output spaces can be expanded to support arbitrary-modality output. For example, Next-GPT (Wu et al., 2023d) and Codi-2 (Tang et al., 2023c) further extends the hybrid output space, while AnyGPT (Zhan et al., 2024a) and UnifiedIO-2 (Lu et al., 2023a) construct unified discrete spaces for all modalities.

### 2.3 Prevalent Representation Paradigms

Considering the representation space structures introduced above, most existing LMMs can be categorized to three types: (1) **Multi-modal understanding models** that rely on Type A input and Type 1 output, these models are mainly designed for understanding tasks that can be fully expressed in language (Dai et al., 2023; Bai et al., 2023b; Lu et al., 2024a; Chen et al., 2023f); (2) **Multi-modal generation models** which comprise of Type A input and Type 2 output, such models excel in generating multi-modal interleaved responses based on the context (Koh et al., 2024; Wu et al., 2023d; Sun et al., 2024a); (3) **Unified multi-modal models** that represent and generate multiple modalities in a unified discrete form (Ge et al., 2023b; Zhan et al., 2024a; Team, 2024). Table 1 and Table 2 in appendix list the design paradigms of contemporary LMMs, grouped according to the aforementioned classification criteria. The alignment architectures discussed in Section 3 are also included.

## 3 Multi-Modal Alignment Architecture

Based on the multi-modal representation spaces introduced in Section 2, the design of LMMs needs to consider how to align representations across different modalities. Mainstream architectures take an LLM-centric paradigm: aligning inputs from all modalities to a **unified multi-modal backbone** for interaction and generating multi-modal responses. To facilitate the unified modeling, additional modules are required, as summarized in Figure 3. We detail the architecture as follows.

### 3.1 Multi-Modal Modeling Backbone

Typically, the backbone is based on a decoder-only architecture composed of multiple transformer blocks (Vaswani et al., 2017). To better understand language, the backbone is primarily initialized with a pre-trained LLM, such as LLaMA (Touvron et al., 2023a,b; Dubey et al., 2024), Vicuna (Chiang et al., 2023), Mistral (Jiang et al., 2023), Qwen (Bai et al., 2023a; Yang et al., 2024), and so on (Bi et al., 2024; Cai et al., 2024; Young et al., 2024). In addition, LMMs for edge devices are usually initialized with smaller language models, such as MobileLLaMA (Kan et al., 2024), Phi (Abdin et al., 2024), etc. (Team et al., 2024a; Hu et al., 2024a). The backbone can also inherit MoE-based language models like Mixtral 8x7B (MistralAITeam, 2023).

Apart from the commonly used architecture mentioned above, some LMMs adopt encoder-decoder backbones (Chung et al., 2024; Chen et al., 2023e; Lu et al., 2023a; Bachmann et al., 2024; Mizrahi et al., 2023). Additionally, native LMMs like Chameleon (Team, 2024) are not initialized with pre-trained LLMs and trained from scratch.

### 3.2 Input-level Alignment

To enable the backbone to process multi-modal information uniformly, it is necessary to align the form and space of inputs across modalities at the input level. Specifically, for Type B input space, since all modalities are represented in a unified discrete token form, input-level alignment can be achieved by directly merging the vocabularies of multiple modalities and learning the token representations through subsequent alignment training (Ge et al., 2023b; Team, 2024; Zhan et al., 2024a).

Regarding Type A hybrid input space, it is required to introduce a connection module to convert inputs from other modalities into a sequential representation that matches the dimension of textual token embeddings. Commonly adopted connection modules are summarized below.

**MLP Based** A typical connection module is multi-layer perceptron (MLP). This module directly aligns the dimension of representations from other modalities with text (Liu et al., 2024b, 2023) by flattening the 2D or 3D features into 1D in a specific order (Maaz et al., 2023; Wu et al., 2023d; Liu et al., 2024a). The advantage of MLP-based module lies in the simplicity and fast convergence during alignment training. However, MLP-based module cannot compress redundant information, which

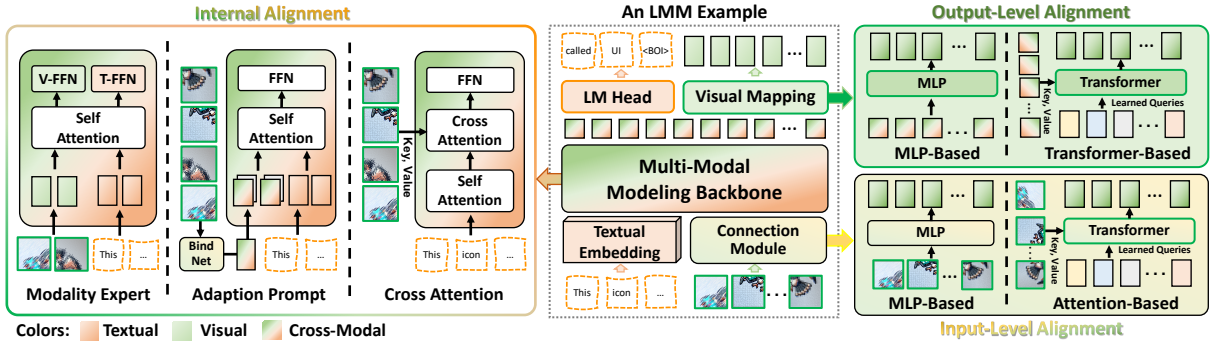


Figure 3: **Summarization of multi-modal alignment modules<sup>1</sup>**: (1) **input-level alignment** to unify the multi-modal inputs into a consistent form and space; (2) **internal alignment** of the backbone for complex cross-modal interactions; and (3) **output-level alignment** to map the outputs to different modality decoders.

could result in excessively long representation sequences (e.g., for high-resolution images). Reducing the computational efficiency requires additional designs to compress the information (Zhu et al., 2023a; Chen et al., 2024c; Dong et al., 2024c).

**Attention Based** Another prevalent connection modules are based on attention mechanisms. This method introduces a fixed number of learnable vectors as queries, which retrieve relevant information from other-modality representations (serving as keys and values) through cross-attention. The output representations of the queries, enriched with information from other modalities, serve as the modality input to the backbone. Representative module architectures include Q-Former (Li et al., 2023b; Dai et al., 2023), abstractor (Ye et al., 2023b,c), resampler (Zeng et al., 2023; Li et al., 2023e), and so on (Bai et al., 2023b; Zhang et al., 2023c). The query-level representations obtained from attention mechanisms effectively compress and aggregate information from other modalities. Additionally, recent works have demonstrated further extensibility, including integrating representations from multiple encoders (Li et al., 2024d; Kar et al., 2024; Tong et al., 2024), incorporating local grounding information (Lu et al., 2023b), and scaling up to an 8B Q-LLaMA (Chen et al., 2023f). However, these modules mainly involve many parameters and typically require additional training (Li et al., 2023b; Lu et al., 2023b). Besides, Yao et al. have found that attention-based modules may result in the loss of important information.

**Others** In addition to the mainstream structures mentioned above, several other connection modules have been proposed. CNN-based modules utilize the inductive bias of convolutional operations to model local information, further combined

with pooling layers, the number of resulted tokens can be effectively reduced (Cha et al., 2024; Chu et al., 2024a; Hong et al., 2024). Adaptive pooling-based modules can compress features using spatial relationships without introducing additional parameters (Yao et al., 2024; Xu et al., 2024a). Furthermore, VL-Mamba explores to use vision selective scanning as connection to integrate representations across different modalities (Qiao et al., 2024).

### 3.3 Internal Alignment

Researchers have explored introducing extra parametric modules within the backbone to further enhance the alignment between modalities. We summarize the commonly adopted methods as follows:

**Cross-Attention Layer** Flamingo (Alayrac et al., 2022) is the first to insert cross-attention layers between the original layers of the backbone, allowing text to perceive information from the visual context. And a tanh gating is introduced to control the degree of modality fusion. Subsequently, such design has been widely adopted by recent LMMs (Gong et al., 2023; Awadalla et al., 2023; IDEFICS, 2023; Chen et al., 2024a), CogAgent (Hong et al., 2023a) further utilizes cross-attention to supplement high-resolution image information. Although effective, densely inserted cross-attention layers bring a large number of additional parameters. Ye et al. improve this by introducing sparsely inserted hyper attention, which significantly reduce extra parameters and facilitate model convergence through parallel self-attention and cross-attention calculation.

**Adaption Prompt** LLaMA-Adapter incorporates visual representations into lightweight learnable adaption prompts and feed the prompts as pre-

<sup>1</sup>The illustration of the input-level and output-level alignment modules is inspired by (Yin et al., 2023).

fix contexts to the backbone (Zhang et al., 2023d). LLaMA-Adapter V2 (Gao et al., 2023) improves this method with an early knowledge fusion strategy. ImageBind-LLM (Han et al., 2023) further extends the prompts to support more modalities.

**Visual Expert** To distinguish between visual and textual modeling, some LMMs introduce visual expert modules to process visual tokens. Specifically, CogVLM (Wang et al., 2023a) adds additional attention and FFN layers to process visual tokens without compromising the original textual modeling capabilities of backbones. mPLUG-Owl2 (Ye et al., 2023c) only introduces modality-specific parameter blocks in the normalization layers and the K and V mapping layers of the attention modules. InternLM-XComposer2 (Dong et al., 2024b), on the other hand, designs a lightweight Partial LoRA module for additional modeling of visual tokens.

### 3.4 Output-level Alignment

Regarding the multi-modal output space described in Section 2.2, both Type 1 and Type 3 are represented in a unified discrete token-based form, multi-modal content can be intuitively generated through a next-token prediction approach with the help of modality-specific de-tokenizer (Zhan et al., 2024a; Lu et al., 2023a; Team, 2024; Ge et al., 2023b).

For the Type 2 hybrid output space, additional mapping modules are required to align the output space of LMM backbones with the input space of corresponding modality generators. Considering image generation, commonly used modules are built on linear projection (Dong et al., 2024a) or the transformer architecture (Koh et al., 2024; Zheng et al., 2024b). Similar to Q-Former, transformer-based modules learn a fixed number of queries to retrieve information from the LMM outputs through cross-attention, serving as the condition input of image diffusion models (Rombach et al., 2022). Next-GPT (Wu et al., 2023d) expands the transformer-based mapping modules to fit the diffusion generators for image, video and audio modalities. Additionally, Emu series (Sun et al., 2024b,a) replace the linear projection with cross-attention in diffusion models to perform dimensional conversion.

In summary, we detail the architectural designs of current large vision-language models (LVLMs) in Table 1. Similarly, such alignment architectures can be extended to more modalities, as presented in Table 2. We kindly refer readers to Appendix B for how to train the constructed models.

## 4 Extension to Embodied Agents

Beyond modality extension, the representation space can be expanded to include various forms of signals in different scenarios, such as embodied environment. In this section, we introduce how to expand LMMs into embodied agents with the intelligence to interact with environments. We will firstly introduce categories of embodied tasks, then delve into how to adapt LMMs to embodied tasks by extending the representation spaces.

### 4.1 Embodied Tasks

Tasks are referred to as “embodied” because the agent needs to interact with a real or virtual environment. Based on the complexity of the interaction actions, we categorize embodied tasks as follows: (1) *Embodied Question Answering (EQA)* (Das et al., 2018; Gordon et al., 2018): In these tasks, the agent is required to answer user questions based on environment exploration. Broadly speaking, we consider such action spaces as discrete vocabularies. (2) *Vision-and-Language Navigation (VLN)* (Anderson et al., 2018; Krantz et al., 2020): These tasks involve navigation based on user instructions. However, these tasks do not require interactions with objects. Therefore, the action space is either discrete directional movements, such as forward, backward, left, and right, or it can involve continuous control parameters, such as speed and direction. (3) *Vision-and-Language Manipulation (VLM)* (Shridhar et al., 2020; Padmakumar et al., 2022; Yenamandra et al., 2023): These tasks require the agent to not only engage in question-answer dialogues with the user, but also navigate the environment and interact with objects based on user instructions. This action space builds upon the action space of VLN tasks by adding object manipulation actions. (4) *Open-World Robot Control (ORC)* (Gupta et al., 2019; Mees et al., 2022; Padalkar et al., 2023): In these tasks, agents are equipped with high-degree-of-freedom robotic arms, capable of performing precise object manipulations, such as grasping and moving objects. The action space for ORC tasks is continuous and determined by the complexity of the robotic arm movements, i.e. represented by a set of continuous values, such as the joint angles or velocities.

### 4.2 Input Extension: Environment

Since embodied agents interact with the environment as the subject, the egocentric observation be-



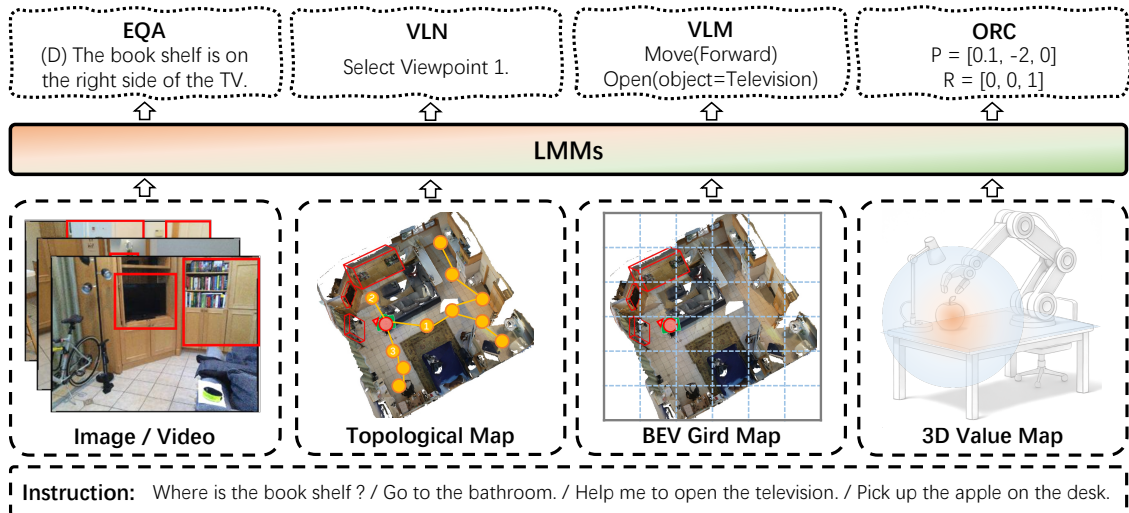


Figure 4: Examples of the input-output space for embodied tasks. Typically, the input includes the user instruction, the current observation (image or video), the environment and the history (optional). We omit the history here, as for different tasks, the content of history vary from pure texts to observation sequences, action sequences and so on.

comes an essential choice (Anderson et al., 2018; Chen et al., 2019; Qi et al., 2020; Padmakumar et al., 2022; Du et al., 2024). Under egocentric observations, the environment is often represented as a local image (Fried et al., 2018; Ahn et al., 2022; Driess et al., 2023) corresponding to the current orientation or by rotating 360 degrees, which could be satisfactory for EQA tasks. However, VLN and VLM tasks require an integrated understanding of observed environments. To obtain a complete picture, the agent must engage in thorough and repeated exploration of the environment (Chaplot et al., 2020a,b). Therefore, the ability to integrate temporal local information and transform it into a long-term global perspective is crucial for embodied agents. Several works utilize topological map (Cartillier et al., 2021, 2024) to record the spatial semantics during navigation, either for obtaining a better visual representation for the environment (Hong et al., 2023b), or for constructing reasoning chains (Zhan et al., 2024b). Others employ bird’s-eye-view grid maps to structure the visited environment (Chen et al., 2023a; Xiong et al., 2023; Wang et al., 2023b). For ORC tasks, a detailed 3D modeling of the environment is essential for executing precise actions with a robotic arm. For example, VoxPoser (Huang et al., 2023b) take the 3D value map derived from interactions between a LLM and a vision-language model to enable exact and efficient object manipulations.

### 4.3 Output Extension: Embodied Action

As stated in Section 4.1, different embodied tasks have distinct embodied action spaces, necessitating

the extensions to model outputs to accommodate the specific demands of each task.

**Discrete Action Space** For embodied tasks of VLN and VLM with discrete action spaces, embodied actions are divided into a fixed set of categories. One line of work, i.e. LLaRP (Szot et al., 2023), utilizes an additional action prediction module to decode discrete actions. Another line of work leverages the powerful language decoding capabilities of LLMs. For example, NavGPT (Zhou et al., 2024b) and NaviLLM (Zheng et al., 2024a) predict actions as plain-text, which are then parsed into specific action commands. This design is simple and effective, yet limits the decoding of complex operations like robotic arm control in ORC tasks. To mitigate this issue, RT-2 (Brohan et al., 2023a) adds special action tokens into the vocabulary. The discrete tokens are then de-tokenized into continuous signals.

**Continuous Action Space** To better adapt to ORC tasks, the extension to continuous actions is necessary. Since the direct outputs of LLMs are discrete tokens, decoding continuous actions typically requires an extra action decoding head. RoboFlamingo (Li et al., 2023d) experiments with different action decoding head architectures (e.g., MLP, RNN, and Transformer) to enable language-conditioned robotic control. Octo (Team et al., 2024b) employs a modular framework, integrating diffusion model-based action policies to predict continuous actions. Unlike RoboFlamingo, the advantage of Octo lies in its ability to flexibly connect different task encoders, observation encoders, and action decoders, making it highly adaptable.

566	<b>Hierarchical Action Space</b>	This separates the	planning, visual question answering and captioning,	616
567		level of action control into high-level task plan-	converting complex environmental perceptions into	617
568		ning and low-level control policies (could be either	multi-step task planning. It integrates the task plans	618
569		discrete or continuous), each handled by separate	with downstream action controller SayCan (Brohan	619
570		modules or models. Specifically, PALM-E (Driess	et al., 2023b) for specific action execution.	620
571		et al., 2023) uses high-level instructions generated		
572		by LVLMs to guide low-level control policies in		
573		executing specific embodied actions.		
574	<b>4.4 Multi-Modal Alignment</b>			
575	<b>Input-level Alignment</b>	To bridge the gap be-	<b>5 Discussion and Outlook</b>	621
576		tween the newly introduced environment repre-		
577		sentation and other modalities, SMNet (Cartillier	<b>Design of representation spaces</b>	622
578		et al., 2021), GridMM (Wang et al., 2023b) and	Constructing	623
579		Trans4Map (Chen et al., 2023a) employ end-to-	multi-modal representation spaces involves hybrid	624
580		end imitation learning, continuously adjusting the	and unified approaches (Section 2 and Section 3),	625
581		model parameters to optimize the updating pro-	each with trade-offs. Hybrid models, which en-	626
582		cesses of allocentric map. However, the obtained	code continuous modality signals into discrete text	627
583		map representations are highly dependent on the	spaces, excel in comprehension but require com-	628
584		UNet and GRU modules nested within the model	plex alignment modules and struggle with genera-	629
585		architecture, lacking the ability to transfer between	tion tasks. Unified discrete models simplify com-	630
586		different language backbones. To address this is-	prehension and generation but face challenges with	631
587		sue, Ego <sup>2</sup> -Map (Hong et al., 2023b) takes a self-	weaker encoders and training stability. Addressing	632
588		supervised contrastive learning strategy, compar-	granularity mismatches between textual and other	633
589		ing egocentric view features with their correspond-	modality tokens is key for future improvement.	634
590		ing semantic maps. Such representations exhibit strong	We kindly refer readers to Appendix C for more	635
591		generalizable capability on various environments.	discussions in this direction.	
592	<b>Output-level Alignment</b>	Adapting the outputs	<b>A promising way towards world models</b>	636
593		to different action spaces is essential for agents	As demonstrated in Section 4, our perspective of rep-	637
594		to understand and execute complex tasks. There	resentation space extension works beyond modal-	638
595		are two major strategies: (1) <i>Direct Alignment</i> :	ities, encompassing any form of information or	639
596		This approach maps instructions directly to exe-	signals. By encoding these into input/output spaces	640
597		cutable actions in an end-to-end manner, as exem-	and aligning them via model architecture and train-	641
598		plified by RoboFlamingo (Li et al., 2023d) and	ing strategies, models can be applied for down-	642
599		Octo (Team et al., 2024b). During training, both	stream tasks. The proposed framework highlights	643
600		RoboFlamingo and Octo collect sequential actions	the potential for models to understand the physical	644
601		covering various scenarios and tasks, enhancing	world. The statement, “ <i>predicting the next token</i>	645
602		the model’s generalization capability during pre-	<i>is to understand the world</i> ”, holds if the defined	646
603		training. They also allow the policy module to	token space has been expanded to cover a sufficient	647
604		be fine-tuned with a small amount of trajectory	amount of information and signals from the world.	648
605		data so as to quickly adapt to new tasks. Besides,		
606		LEO (Huang et al., 2023a) adopts a two-stage train-	<b>6 Conclusion</b>	649
607		ing process involving pre-training for 3D vision-		
608		language alignment and fine-tuning on 3D vision-	In this paper, we summarize the current methods of	650
609		language-action instructions, enhancing the agent’s	LMM construction from the perspective of repre-	651
610		adaptability to different action spaces. (2) <i>Indi-</i>	sentation space extension. We further break down	652
611		<i>rect Alignment</i> : This method breaks down user	and provide detailed discussion of the key research	653
612		instructions into language plans that can be under-	problems in the construction process, including the	654
613		stood by downstream models, with representative	structure of multi-modal input and output repre-	655
614		works as PALM-E (Driess et al., 2023). PALM-E	sentation spaces and multi-modal representation	656
615		pre-trains on large datasets of robotic manipulation	alignment frameworks. Our summarization frame-	657
			work is not only straightforward but also effectively	658
			encapsulates the mainstream approaches while of-	659
			fering potential for future extensions. This paper	660
			will continue to be updated, and we hope it can	661
			provide an intuitive and comprehensive overview	662
			for related researchers and inspire future work.	663



## 7 Limitations

Although our analysis from the perspective of representation space extension is general, this paper does not delve deeply into the evaluation of existing LMMs. Notably, evaluation tasks and datasets can be systematically categorized based on the input and output representation spaces. For example, VQA tasks that only require Type 1 outputs (see Figure 2) could be considered “understanding” tasks, while image editing tasks that require Type 2 or Type 3 outputs could be categorised as “generation” tasks. This aspect remains an open question and we reserve it for future investigation.

## References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.

Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. 2022. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*.

AI@Meta. 2024. [Llama 3 model card](#).

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *NIPS*, 35:23716–23736.

Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton Van Den Hengel. 2018. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3674–3683.

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

Roman Bachmann, Oğuzhan Fatih Kar, David Mizrahi, Ali Garjani, Mingfei Gao, David Griffiths, Jiaming Hu, Afshin Dehghan, and Amir Zamir. 2024. [4m-21: An any-to-any vision model for tens of tasks and modalities](#). *Preprint*, arXiv:2406.09406.

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

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023b. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*.

Tianyi Bai, Hao Liang, Binwang Wan, Ling Yang, Bozhou Li, Yifan Wang, Bin Cui, Conghui He, Binhang Yuan, and Wentao Zhang. 2024. A survey of multimodal large language model from a data-centric perspective. *arXiv preprint arXiv:2405.16640*.

Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşlılar. 2023. [Fuyu-8b](#).

Gedas Bertasius, Heng Wang, and Lorenzo Torresani. 2021. Is space-time attention all you need for video understanding? In *ICML*, volume 2, page 4.

Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. 2024. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*.

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. 2023a. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*.

Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. 2023b. Do as i can, not as i say: Grounding language in robotic affordances. In *Conference on robot learning*, pages 287–318. PMLR.

Davide Caffagni, Federico Cocchi, Luca Barsellotti, Nicholas Moratelli, Sara Sarto, Lorenzo Baraldi, Marcella Cornia, and Rita Cucchiara. 2024. The (r) evolution of multimodal large language models: A survey. *arXiv preprint arXiv:2402.12451*.

Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, et al. 2024. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*.

Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. 2021. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660.

Vincent Cartillier, Neha Jain, and Irfan Essa. 2024. 3d semantic mapnet: Building maps for multi-object re-identification in 3d. *arXiv preprint arXiv:2403.13190*.

771	Vincent Cartillier, Zhile Ren, Neha Jain, Stefan Lee, Irfan Essa, and Dhruv Batra. 2021. Semantic mapnet: Building allocentric semantic maps and representations from egocentric views. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pages 964–972.	827
772		828
773		829
774		830
775		831
776		832
777	Junbum Cha, Wooyoung Kang, Jonghwan Mun, and Byungseok Roh. 2024. Honeybee: Locality-enhanced projector for multimodal llm. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 13817–13827.	833
778		834
779		835
780		836
781		837
782	Devendra Singh Chaplot, Dhiraj Gandhi, Saurabh Gupta, Abhinav Gupta, and Ruslan Salakhutdinov. 2020a. Learning to explore using active neural slam. <i>arXiv preprint arXiv:2004.05155</i> .	838
783		839
784		840
785		841
786	Devendra Singh Chaplot, Ruslan Salakhutdinov, Abhinav Gupta, and Saurabh Gupta. 2020b. Neural topological slam for visual navigation. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 12875–12884.	842
787		843
788		844
789		845
790		846
791	Chang Chen, Jiaming Zhang, Kailun Yang, Kunyu Peng, and Rainer Stiefelwagen. 2023a. Trans4map: Revisiting holistic bird’s-eye-view mapping from egocentric images to allocentric semantics with vision transformers. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pages 4013–4022.	847
792		848
793		849
794		850
795		851
796		852
797		853
798	Howard Chen, Alane Suhr, Dipendra Misra, Noah Snively, and Yoav Artzi. 2019. Touchdown: Natural language navigation and spatial reasoning in visual street environments. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 12538–12547.	854
799		855
800		856
801		857
802		858
803		859
804	Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. 2023b. Minigt-v2: large language model as a unified interface for vision-language multi-task learning. <i>arXiv:2310.09478</i> .	860
805		861
806		862
807		863
808		864
809		865
810	Kaibing Chen, Dong Shen, Hanwen Zhong, Huasong Zhong, Kui Xia, Di Xu, Wei Yuan, Yifei Hu, Bin Wen, Tianke Zhang, et al. 2024a. Evlm: An efficient vision-language model for visual understanding. <i>arXiv preprint arXiv:2407.14177</i> .	866
811		867
812		868
813		869
814		870
815	Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. 2023c. Shikra: Unleashing multimodal llm’s referential dialogue magic. <i>arXiv:2306.15195</i> .	871
816		872
817		873
818		874
819	Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023d. Sharegpt4v: Improving large multi-modal models with better captions. <i>arXiv:2311.12793</i> .	875
820		876
821		877
822		878
823	Shaoxiang Chen, Zequn Jie, and Lin Ma. 2024b. Llavamole: Sparse mixture of lora experts for mitigating data conflicts in instruction finetuning mllms. <i>arXiv preprint arXiv:2401.16160</i> .	879
824		880
825		881
826		882
	Xi Chen, Josip Djolonga, Piotr Padlewski, Basil Mustafa, Soravit Changpinyo, Jialin Wu, Carlos Riquelme Ruiz, Sebastian Goodman, Xiao Wang, Yi Tay, Siamak Shakeri, Mostafa Dehghani, Daniel Salz, Mario Lucic, Michael Tschannen, Arsha Nagrani, Hexiang Hu, Mandar Joshi, Bo Pang, Ceslee Montgomery, Paulina Pietrzyk, Marvin Ritter, AJ Piergiovanni, Matthias Minderer, Filip Pavetic, Austin Waters, Gang Li, Ibrahim Alabdulmohsin, Lucas Beyer, Julien Amelot, Kenton Lee, Andreas Peter Steiner, Yang Li, Daniel Keysers, Anurag Arnab, Yuanzhong Xu, Keran Rong, Alexander Kolesnikov, Mojtaba Seyedhosseini, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. 2023e. Pali-x: On scaling up a multilingual vision and language model. <i>Preprint</i> , arXiv:2305.18565.	883
	Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. 2024c. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. <i>arXiv preprint arXiv:2404.16821</i> .	884
	Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Zhong Muyan, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 2023f. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. <i>arXiv preprint arXiv:2312.14238</i> .	885
	Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality.	886
	Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, et al. 2023a. Mobilevlm: A fast, reproducible and strong vision language assistant for mobile devices. <i>arXiv preprint arXiv:2312.16886</i> .	887
	Xiangxiang Chu, Limeng Qiao, Xinyu Zhang, Shuang Xu, Fei Wei, Yang Yang, Xiaofei Sun, Yiming Hu, Xinyang Lin, Bo Zhang, et al. 2024a. Mobilevlm v2: Faster and stronger baseline for vision language model. <i>arXiv preprint arXiv:2402.03766</i> .	888
	Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, et al. 2024b. Qwen2-audio technical report. <i>arXiv preprint arXiv:2407.10759</i> .	889
	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> .	890
	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai,	891
		892
		893
		894

885	Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2024. <a href="#">Scaling instruction-finetuned language models</a> . <i>Journal of Machine Learning Research</i> , 25(70):1–53.	943
886		944
887		945
888		946
889		947
890		
891		948
892		949
893		950
894	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.	951
895		952
896		953
897		954
898		
899		955
900		956
901	Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. <a href="#">Instructblip: Towards general-purpose vision-language models with instruction tuning</a> . <i>Preprint</i> , arXiv:2305.06500.	957
902		958
903		959
904		
905		960
906		961
907		962
908	Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. 2018. Embodied question answering. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 1–10.	963
909		964
910		
911		965
912	Nilaksh Das, Saket Dingliwal, Srikanth Ronanki, Rohit Paturi, David Huang, Prashant Mathur, Jie Yuan, Dhanush Bekal, Xing Niu, Sai Muralidhar Jayanthi, et al. 2024. <a href="#">Speechverse: A large-scale generalizable audio language model</a> . <i>arXiv preprint arXiv:2405.08295</i> .	966
913		967
914		968
915		969
916		
917		970
918	Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu Zhou, Haoran Wei, Xiangwen Kong, Xiangyu Zhang, Kaisheng Ma, and Li Yi. 2024a. <a href="#">Dreamllm: Synergistic multimodal comprehension and creation</a> . <i>Preprint</i> , arXiv:2309.11499.	971
919		972
920		973
921		974
922		975
923		
924	Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang Zhang, Haodong Duan, Maosong Cao, et al. 2024b. <a href="#">Internlm-xcomposer2: Mastering free-form text-image composition and comprehension in vision-language large model</a> . <i>arXiv preprint arXiv:2401.16420</i> .	976
925		977
926		978
927		979
928		980
929		
930		981
931	Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Songyang Zhang, Haodong Duan, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Zhe Chen, Xinyue Zhang, Wei Li, Jingwen Li, Wenhai Wang, Kai Chen, Conghui He, Xingcheng Zhang, Jifeng Dai, Yu Qiao, Dahua Lin, and Jiaqi Wang. 2024c. <a href="#">Internlm-xcomposer2-4khd: A pioneering large vision-language model handling resolutions from 336 pixels to 4k hd</a> . <i>Preprint</i> , arXiv:2404.06512.	982
932		983
933		984
934		985
935		
936		986
937		987
938		988
939		989
940		
941	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,	990
942		991
		992
		993
		994
		995
		996
		997
		998



999	v2: Parameter-efficient visual instruction model. <i>arXiv:2304.15010</i> .		1055	
1000			1056	
1001	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> .		1057	
1002			1058	
1003				
1004				
1005				
1006				
1007	Chunjiang Ge, Sijie Cheng, Ziming Wang, Jiale Yuan, Yuan Gao, Jun Song, Shiji Song, Gao Huang, and Bo Zheng. 2024. Convllava: Hierarchical backbones as visual encoder for large multimodal models. <i>arXiv preprint arXiv:2405.15738</i> .			
1008				
1009				
1010				
1011				
1012	Yuying Ge, Yixiao Ge, Ziyun Zeng, Xintao Wang, and Ying Shan. 2023a. <a href="#">Planting a seed of vision in large language model</a> . <i>Preprint</i> , arXiv:2307.08041.			
1013				
1014				
1015	Yuying Ge, Sijie Zhao, Ziyun Zeng, Yixiao Ge, Chen Li, Xintao Wang, and Ying Shan. 2023b. Making llama see and draw with seed tokenizer. <i>arXiv preprint arXiv:2310.01218</i> .			
1016				
1017				
1018				
1019	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 (CVPR)</i> , pages 15180–15190.			
1020				
1021				
1022				
1023				
1024				
1025	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:2305.04790</i> .			
1026				
1027				
1028				
1029				
1030	Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. 2018. Iqa: Visual question answering in interactive environments. In <i>Proceedings of the IEEE conference on computer vision and pattern recogni- tion</i> , pages 4089–4098.			
1031				
1032				
1033				
1034				
1035				
1036	Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. 2019. Relay pol- icy learning: Solving long-horizon tasks via imi- tation and reinforcement learning. <i>arXiv preprint arXiv:1910.11956</i> .			
1037				
1038				
1039				
1040				
1041	Jiaming Han, Renrui Zhang, Wenqi Shao, Peng Gao, Peng Xu, Han Xiao, Kaipeng Zhang, Chris Liu, Song Wen, Ziyu Guo, et al. 2023. Imagebind-llm: Multi- modality instruction tuning. <i>arXiv:2309.03905</i> .			
1042				
1043				
1044				
1045	Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsu- pervised visual representation learning. In <i>Proceeed- ings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 9729–9738.			
1046				
1047				
1048				
1049				
1050	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.			
1051				
1052				
1053				
1054				
		Muyang He, Yexin Liu, Boya Wu, Jianhao Yuan, Yuezhe Wang, Tiejun Huang, and Bo Zhao. 2024. Efficient multimodal learning from data-centric perspective. <i>arXiv preprint arXiv:2402.11530</i> .		1059
			1060	
			1061	
			1062	
			1063	
		Wenyi Hong, Weihang Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng, Shiyu Huang, Junhui Ji, Zhao Xue, et al. 2024. Cogvlm2: Visual language models for image and video understanding. <i>arXiv preprint arXiv:2408.16500</i> .		1064
			1065	
			1066	
			1067	
			1068	
		Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. 2023a. Cogagent: A visual language model for gui agents. <i>arXiv preprint arXiv:2312.08914</i> .		1069
			1070	
			1071	
			1072	
			1073	
			1074	
		Yicong Hong, Yang Zhou, Ruiyi Zhang, Franck Der- noncourt, Trung Bui, Stephen Gould, and Hao Tan. 2023b. Learning navigational visual representations with semantic map supervision. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 3055–3067.		1075
			1076	
			1077	
			1078	
			1079	
			1080	
		Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdel- rahman 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.		1081
			1082	
			1083	
			1084	
			1085	
			1086	
		Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxi- ang Huang, Weilin Zhao, et al. 2024a. Minicpm: Unveiling the potential of small language models with scalable training strategies. <i>arXiv preprint arXiv:2404.06395</i> .		1087
			1088	
			1089	
			1090	
			1091	
		Shujie Hu, Long Zhou, Shujie Liu, Sanyuan Chen, Hongkun Hao, Jing Pan, Xunying Liu, Jinyu Li, Sunit Sivasankaran, Linqun Liu, et al. 2024b. Wavllm: Towards robust and adaptive speech large language model. <i>arXiv preprint arXiv:2404.00656</i> .		1092
			1093	
			1094	
			1095	
		Wenbo Hu, Yifan Xu, Y Li, W Li, Z Chen, and Z Tu. 2023. Bliva: A simple multimodal llm for better handling of text-rich visual questions. <i>arXiv:2308.09936</i> .		1096
			1097	
			1098	
			1099	
			1100	
		Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. 2023a. An embodied generalist agent in 3d world. <i>arXiv preprint arXiv:2311.12871</i> .		1101
			1102	
			1103	
		Jiaxing Huang and Jingyi Zhang. 2024. A survey on evaluation of multimodal large language models. <i>arXiv preprint arXiv:2408.15769</i> .		1104
			1105	
			1106	
			1107	
			1108	
		Wenlong Huang, Chen Wang, Ruohan Zhang, Yun- zhu Li, Jiajun Wu, and Li Fei-Fei. 2023b. Vox- poser: Composable 3d value maps for robotic ma- nipulation with language models. <i>arXiv preprint arXiv:2307.05973</i> .		

1109	IDEFICS. 2023. Introducing idefics: An open reproduction of state-of-the-art visual language model. <a href="https://huggingface.co/blog/idefics">https://huggingface.co/blog/idefics</a> .	1165
1110		1166
1111		1167
1112	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	1168
1113		1169
1114		1170
1115		1171
1116		1172
1117	Yang Jin, Kun Xu, Kun Xu, Liwei Chen, Chao Liao, Jianchao Tan, Yadong Mu, et al. 2023. Unified language-vision pretraining in llm with dynamic discrete visual tokenization. <i>arXiv preprint arXiv:2309.04669</i> .	1173
1118		1174
1119		1175
1120		1176
1121		1177
1122	Khen Bo Kan, Hyunsu Mun, Guohong Cao, and Youngseok Lee. 2024. Mobile-llama: Instruction fine-tuning open-source llm for network analysis in 5g networks. <i>IEEE Network</i> .	1178
1123		1179
1124		1180
1125		1181
1126	Oğuzhan Fatih Kar, Alessio Tonioni, Petra Poklukar, Achin Kulshrestha, Amir Zamir, and Federico Tombari. 2024. Brave: Broadening the visual encoding of vision-language models. <i>arXiv preprint arXiv:2404.07204</i> .	1182
1127		1183
1128		1184
1129		1185
1130		1186
1131	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. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 4015–4026.	1187
1132		1188
1133		1189
1134		1190
1135		1191
1136		1192
1137	Jing Yu Koh, Daniel Fried, and Russ R Salakhutdinov. 2024. Generating images with multimodal language models. <i>Advances in Neural Information Processing Systems</i> , 36.	1193
1138		1194
1139		1195
1140		1196
1141	Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. 2020. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In <i>Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII 16</i> , pages 104–120. Springer.	1197
1142		1198
1143		1199
1144		1200
1145		1201
1146		1202
1147		1203
1148	Taku Kudo. 2018. Subword regularization: Improving neural network translation models with multiple subword candidates. <i>arXiv preprint arXiv:1804.10959</i> .	1204
1149		1205
1150		1206
1151	Siddique Latif, Moazzam Shoukat, Fahad Shamshad, Muhammad Usama, Yi Ren, Heriberto Cuayáhuitl, Wenwu Wang, Xulong Zhang, Roberto Togneri, Erik Cambria, et al. 2023. Sparks of large audio models: A survey and outlook. <i>arXiv preprint arXiv:2308.12792</i> .	1207
1152		1208
1153		1209
1154		1210
1155		1211
1156		1212
1157	Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. 2024a. Building and better understanding vision-language models: insights and future directions. <i>arXiv preprint arXiv:2408.12637</i> .	1213
1158		1214
1159		1215
1160		1216
1161	Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. 2024b. What matters when building vision-language models? <i>arXiv preprint arXiv:2405.02246</i> .	1217
1162		1218
1163		1219
1164		1220
	Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023a. <i>Otter: A multi-modal model with in-context instruction tuning</i> . Preprint, arXiv:2305.03726.	1221
		1222
	Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024a. <i>Llava-onevision: Easy visual task transfer</i> . Preprint, arXiv:2408.03326.	1223
		1224
	Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. 2024b. <i>Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models</i> . <i>arXiv preprint arXiv:2407.07895</i> .	1225
		1226
	Jiachen Li, Xinyao Wang, Sijie Zhu, Chia-Wen Kuo, Lu Xu, Fan Chen, Jitesh Jain, Humphrey Shi, and Longyin Wen. 2024c. <i>Cumo: Scaling multimodal llm with co-upcycled mixture-of-experts</i> . <i>arXiv preprint arXiv:2405.05949</i> .	1227
		1228
	Jian Li and Weiheng Lu. 2024. A survey on benchmarks of multimodal large language models. <i>arXiv preprint arXiv:2408.08632</i> .	1229
		1230
	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. <i>Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models</i> . <i>arXiv:2301.12597</i> .	1231
		1232
	KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023c. <i>Videochat: Chat-centric video understanding</i> . <i>arXiv:2305.06355</i> .	1233
		1234
	Xinghang Li, Minghuan Liu, Hanbo Zhang, Cunjun Yu, Jie Xu, Hongtao Wu, Chilam Cheang, Ya Jing, Weinan Zhang, Huaping Liu, et al. 2023d. <i>Vision-language foundation models as effective robot imitators</i> . <i>arXiv preprint arXiv:2311.01378</i> .	1235
		1236
	Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. 2024d. <i>Mini-gemini: Mining the potential of multi-modality vision language models</i> . <i>arXiv preprint arXiv:2403.18814</i> .	1237
		1238
	Yunxin Li, Shenyuan Jiang, Baotian Hu, Longyue Wang, Wanqi Zhong, Wenhan Luo, Lin Ma, and Min Zhang. 2024e. <i>Uni-moe: Scaling unified multimodal llms with mixture of experts</i> . <i>arXiv preprint arXiv:2405.11273</i> .	1239
		1240
	Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. 2023e. <i>Monkey: Image resolution and text label are important things for large multi-modal models</i> . <i>arXiv preprint arXiv:2311.06607</i> .	1241
		1242
	Bin Lin, Zhenyu Tang, Yang Ye, Jiayi Cui, Bin Zhu, Peng Jin, Junwu Zhang, Munan Ning, and Li Yuan. 2024a. <i>Moe-llava: Mixture of experts for large vision-language models</i> . <i>arXiv preprint arXiv:2401.15947</i> .	1243
		1244
		1245
		1246

1219	Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad Shoeybi, and Song Han. 2024b. <a href="#">Vila: On pre-training for visual language models</a> . <i>Preprint</i> , arXiv:2312.07533.	fine-grained language-vision alignment and comprehension via semantic-aware visual objects. <i>arXiv preprint arXiv:2312.05278</i> .	1273 1274 1275
1224	Xi Victoria Lin, Akshat Shrivastava, Liang Luo, Srinivasan Iyer, Mike Lewis, Gargi Gosh, Luke Zettlemoyer, and Armen Aghajanyan. 2024c. Moma: Efficient early-fusion pre-training with mixture of modality-aware experts. <i>arXiv preprint arXiv:2407.21770</i> .	Shiyin Lu, Yang Li, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, and Han-Jia Ye. 2024b. Ovis: Structural embedding alignment for multimodal large language model. <i>arXiv preprint arXiv:2405.20797</i> .	1276 1277 1278 1279
1230	Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi Shao, Keqin Chen, et al. 2023. Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. <i>arXiv preprint arXiv:2311.07575</i> .	Gen Luo, Yiyi Zhou, Tianhe Ren, Shengxin Chen, Xiaoshuai Sun, and Rongrong Ji. 2023a. Cheap and quick: Efficient vision-language instruction tuning for large language models. <i>arXiv:2305.15023</i> .	1280 1281 1282 1283
1236	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023. Improved baselines with visual instruction tuning. <i>arXiv:2310.03744</i> .	Ruipu Luo, Ziwang Zhao, Min Yang, Junwei Dong, Minghui Qiu, Pengcheng Lu, Tao Wang, and Zhongyu Wei. 2023b. Valley: Video assistant with large language model enhanced ability. <i>arXiv preprint arXiv:2306.07207</i> .	1284 1285 1286 1287 1288
1239	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024a. Llava-next: Improved reasoning, ocr, and world knowledge.	Ziyang Ma, Guanrou Yang, Yifan Yang, Zhifu Gao, Jiaming Wang, Zhihao Du, Fan Yu, Qian Chen, Siqi Zheng, Shiliang Zhang, et al. 2024. An embarrassingly simple approach for llm with strong asr capacity. <i>arXiv preprint arXiv:2402.08846</i> .	1289 1290 1291 1292 1293
1242	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024b. Visual instruction tuning. <i>Advances in neural information processing systems</i> , 36.	Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2023. Video-chatgpt: Towards detailed video understanding via large vision and language models. <i>arXiv preprint arXiv:2306.05424</i> .	1294 1295 1296 1297 1298
1245	Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 10012–10022.	Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruvi Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. 2024. Mm1: Methods, analysis & insights from multimodal llm pre-training. <i>arXiv preprint arXiv:2403.09611</i> .	1299 1300 1301 1302 1303 1304
1251	Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. 2022a. Video swin transformer. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 3202–3211.	Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. 2022. Calvin: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks. <i>IEEE Robotics and Automation Letters</i> , 7(3):7327–7334.	1305 1306 1307 1308 1309
1256	Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. 2022b. A convnet for the 2020s. arxiv e-prints.	MistralAITeam. 2023. Mixtral of experts a high quality sparse mixture-of-experts. [EB/OL]. <a href="https://mistral.ai/news/mixtral-of-experts/">https://mistral.ai/news/mixtral-of-experts/</a> Accessed December 11, 2023.	1310 1311 1312 1313
1259	Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, et al. 2024a. Deepseek-vl: Towards real-world vision-language understanding. <i>arXiv:2403.05525</i> .	David Mizrahi, Roman Bachmann, Oğuzhan Fatih Kar, Teresa Yeo, Mingfei Gao, Afshin Dehghan, and Amir Zamir. 2023. 4m: Massively multimodal masked modeling. <i>Preprint</i> , arXiv:2312.06647.	1314 1315 1316 1317
1264	Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. 2023a. Unified-io 2: Scaling autoregressive multimodal models with vision, language, audio, and action. <i>Preprint</i> , arXiv:2312.17172.	OpenAI. 2023. <a href="#">Chatgpt (august 3 version)</a> .	1318
1270	Junyu Lu, Ruyi Gan, Dixiang Zhang, Xiaojun Wu, Ziwei Wu, Renliang Sun, Jiaying Zhang, Pingjian Zhang, and Yan Song. 2023b. Lyrics: Boosting	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> .	1319 1320 1321 1322 1323 1324



1325	Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex	Zettlemoyer, and Dieter Fox. 2020. Alfred: A bench-	1380
1326	Bewley, Alex Herzog, Alex Irpan, Alexander Khaz-	mark for interpreting grounded instructions for ev-	1381
1327	atsky, Anant Rai, Anikait Singh, Anthony Bro-	everyday tasks. In <i>Proceedings of the IEEE/CVF con-</i>	1382
1328	han, et al. 2023. Open x-embodiment: Robotic	<i>ference on computer vision and pattern recognition</i> ,	1383
1329	learning datasets and rt-x models. <i>arXiv preprint</i>	pages 10740–10749.	1384
1330	<i>arXiv:2310.08864</i> .		
1331	Aishwarya Padmakumar, Jesse Thomason, Ayush Shri-	Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang,	1385
1332	vastava, Patrick Lange, Anjali Narayan-Chen, Span-	and Deng Cai. 2023. Pandagpt: One model to	1386
1333	dana Gella, Robinson Piramuthu, Gokhan Tur, and	instruction-follow them all. <i>arXiv:2305.16355</i> .	1387
1334	Dilek Hakkani-Tur. 2022. Teach: Task-driven em-	Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang,	1388
1335	odied agents that chat. In <i>Proceedings of the AAAI</i>	Qiyang Yu, Zhengxiong Luo, Yueze Wang, Yongming	1389
1336	<i>Conference on Artificial Intelligence</i> , volume 36,	Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang.	1390
1337	pages 2017–2025.	2024a. <b>Generative multimodal models are in-context</b>	1391
1338	Xichen Pan, Li Dong, Shaohan Huang, Zhiliang Peng,	<b>learners</b> . <i>Preprint</i> , arXiv:2312.13286.	1392
1339	Wenhu Chen, and Furu Wei. 2024. <b>Kosmos-g: Gen-</b>	Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and	1393
1340	<b>erating images in context with multimodal large lan-</b>	Yue Cao. 2023. Eva-clip: Improved training tech-	1394
1341	<b>guage models</b> . <i>Preprint</i> , arXiv:2310.02992.	niques for clip at scale. <i>arXiv:2303.15389</i> .	1395
1342	Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang,	Quan Sun, Qiyang Yu, Yufeng Cui, Fan Zhang,	1396
1343	William Yang Wang, Chunhua Shen, and Anton	Xiaosong Zhang, Yueze Wang, Hongcheng Gao,	1397
1344	van den Hengel. 2020. Reverie: Remote embodied	Jingjing Liu, Tiejun Huang, and Xinlong Wang.	1398
1345	visual referring expression in real indoor environ-	2024b. <b>Emu: Generative pretraining in multimodal-</b>	1399
1346	ments. In <i>Proceedings of the IEEE/CVF Conference</i>	<b>ity</b> . <i>Preprint</i> , arXiv:2307.05222.	1400
1347	<i>on Computer Vision and Pattern Recognition</i> , pages	Andrew Szot, Max Schwarzer, Harsh Agrawal, Bog-	1401
1348	9982–9991.	dan Mazouze, Rin Metcalf, Walter Talbott, Natalie	1402
1349	Yanyuan Qiao, Zheng Yu, Longteng Guo, Sihan	Mackraz, R Devon Hjelm, and Alexander T Toshev.	1403
1350	Chen, Zijia Zhao, Mingzhen Sun, Qi Wu, and	2023. Large language models as generalizable poli-	1404
1351	Jing Liu. 2024. V1-mamba: Exploring state space	cies for embodied tasks. In <i>The Twelfth International</i>	1405
1352	models for multimodal learning. <i>arXiv preprint</i>	<i>Conference on Learning Representations</i> .	1406
1353	<i>arXiv:2403.13600</i> .	Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao	1407
1354	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and Chao	1408
1355	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	Zhang. 2023a. Salmonn: Towards generic hearing	1409
1356	try, Amanda Askell, Pamela Mishkin, Jack Clark,	abilities for large language models. <i>arXiv preprint</i>	1410
1357	et al. 2021. Learning transferable visual models from	<i>arXiv:2310.13289</i> .	1411
1358	natural language supervision. In <i>ICML</i> , pages 8748–	Yunlong Tang, Jing Bi, Siting Xu, Luchuan Song, Susan	1412
1359	8763. PMLR.	Liang, Teng Wang, Daoan Zhang, Jie An, Jingyang	1413
1360	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	Lin, Rongyi Zhu, et al. 2023b. Video understanding	1414
1361	Dario Amodei, Ilya Sutskever, et al. 2019. Language	with large language models: A survey. <i>arXiv preprint</i>	1415
1362	models are unsupervised multitask learners. <i>OpenAI</i>	<i>arXiv:2312.17432</i> .	1416
1363	<i>blog</i> , 1(8):9.	Zineng Tang, Ziyi Yang, Mahmoud Khademi, Yang Liu,	1417
1364	Robin Rombach, Andreas Blattmann, Dominik Lorenz,	Chenguang Zhu, and Mohit Bansal. 2023c. <b>Codi-</b>	1418
1365	Patrick Esser, and Björn Ommer. 2022. High-	<b>2: In-context, interleaved, and interactive any-to-any</b>	1419
1366	resolution image synthesis with latent diffusion mod-	<b>generation</b> . <i>Preprint</i> , arXiv:2311.18775.	1420
1367	els. In <i>Proceedings of the IEEE/CVF Conference on</i>	Chameleon Team. 2024. Chameleon: Mixed-modal	1421
1368	<i>Computer Vision and Pattern Recognition (CVPR)</i> ,	early-fusion foundation models. <i>arXiv preprint</i>	1422
1369	pages 10684–10695.	<i>arXiv:2405.09818</i> .	1423
1370	Rico Sennrich, Barry Haddow, and Alexandra Birch.	Gemma Team, Thomas Mesnard, Cassidy Hardin,	1424
1371	2015. Neural machine translation of rare words with	Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,	1425
1372	subword units. <i>arXiv preprint arXiv:1508.07909</i> .	Laurent Sifre, Morgane Rivi�re, Mihir Sanjay Kale,	1426
1373	Zhenwei Shao, Zhou Yu, Jun Yu, Xuecheng Ouyang,	Juliette Love, et al. 2024a. Gemma: Open models	1427
1374	Lihao Zheng, Zhenbiao Gai, Mingyang Wang, and	based on gemini research and technology. <i>arXiv</i>	1428
1375	Jiajun Ding. 2024. Imp: Highly capable large mul-	<i>preprint arXiv:2403.08295</i> .	1429
1376	timodal models for mobile devices. <i>arXiv preprint</i>	Octo Model Team, Dibya Ghosh, Homer Walke, Karl	1430
1377	<i>arXiv:2405.12107</i> .	Pertsch, Kevin Black, Oier Mees, Sudeep Dasari,	1431
1378	Mohit Shridhar, Jesse Thomason, Daniel Gordon,	Joey Hejna, Tobias Kreiman, Charles Xu, et al. 2024b.	1432
1379	Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke	Octo: An open-source generalist robot policy. <i>arXiv</i>	1433
		<i>preprint arXiv:2405.12213</i> .	1434

1435	Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. 2024. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. <i>arXiv preprint arXiv:2406.16860</i> .	1492
1436		1493
1437		1494
1438		1495
1439		1496
1440		
1441	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:2302.13971</i> .	1497
1442		1498
1443		1499
1444		
1445		
1446	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruiti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv:2307.09288</i> .	1500
1447		1501
1448		1502
1449		1503
1450		1504
1451		1505
1452	Aaron Van Den Oord, Oriol Vinyals, et al. 2017. Neural discrete representation learning. <i>Advances in neural information processing systems</i> , 30.	1506
1453		1507
1454		1508
1455	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. <b>Attention is all you need</b> . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	1509
1456		1510
1457		1511
1458		1512
1459		1513
1460	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .	1514
1461		1515
1462		1516
1463		1517
1464		1518
1465	Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. 2023a. Cogvlm: Visual expert for pretrained language models. <i>arXiv preprint arXiv:2311.03079</i> .	1519
1466		1520
1467		1521
1468		1522
1469		1523
1470		1524
1471	Zihan Wang, Xiangyang Li, Jiahao Yang, Yeqi Liu, and Shuqiang Jiang. 2023b. Gridmm: Grid memory map for vision-and-language navigation. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 15625–15636.	1525
1472		1526
1473		1527
1474		1528
1475	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> .	1529
1476		1530
1477		1531
1478		1532
1479		1533
1480	Jialin Wu, Xia Hu, Yaqing Wang, Bo Pang, and Radu Soricut. 2024a. Omni-smola: Boosting generalist multimodal models with soft mixture of low-rank experts. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 14205–14215.	1534
1481		1535
1482		1536
1483		1537
1484		1538
1485	Jian Wu, Yashesh Gaur, Zhuo Chen, Long Zhou, Yimeng Zhu, Tianrui Wang, Jinyu Li, Shujie Liu, Bo Ren, Linqun Liu, et al. 2023b. On decoder-only architecture for speech-to-text and large language model integration. In <i>2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 1–8. IEEE.	1539
1486		1540
1487		1541
1488		1542
1489		1543
1490		1544
1491		1545
	Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and S Yu Philip. 2023c. Multimodal large language models: A survey. In <i>2023 IEEE International Conference on Big Data (BigData)</i> , pages 2247–2256. IEEE.	1492
		1493
		1494
		1495
		1496
	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> .	1497
		1498
		1499
	Yecheng Wu, Zhuoyang Zhang, Junyu Chen, Haotian Tang, Dacheng Li, Yunhao Fang, Ligeng Zhu, Enze Xie, Hongxu Yin, Li Yi, et al. 2024b. Vila-u: a unified foundation model integrating visual understanding and generation. <i>arXiv preprint arXiv:2409.04429</i> .	1500
		1501
		1502
		1503
		1504
		1505
	Yonghui Wu. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. <i>arXiv preprint arXiv:1609.08144</i> .	1506
		1507
		1508
	Hanguang Xiao, Feizhong Zhou, Xingyue Liu, Tianqi Liu, Zhipeng Li, Xin Liu, and Xiaoxuan Huang. 2024. A comprehensive survey of large language models and multimodal large language models in medicine. <i>arXiv preprint arXiv:2405.08603</i> .	1509
		1510
		1511
		1512
		1513
	Xuan Xiong, Yicheng Liu, Tianyuan Yuan, Yue Wang, Yilun Wang, and Hang Zhao. 2023. Neural map prior for autonomous driving. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 17535–17544.	1514
		1515
		1516
		1517
		1518
	Haiyang Xu, Qinghao Ye, Xuan Wu, Ming Yan, Yuan Miao, Jiabo Ye, Guohai Xu, Anwen Hu, Yaya Shi, Guangwei Xu, et al. 2023. Youku-mplug: A 10 million large-scale chinese video-language dataset for pre-training and benchmarks. <i>arXiv preprint arXiv:2306.04362</i> .	1519
		1520
		1521
		1522
		1523
		1524
	Lin Xu, Yilin Zhao, Daquan Zhou, Zhijie Lin, See Kiong Ng, and Jiashi Feng. 2024a. Pllava: Parameter-free llava extension from images to videos for video dense captioning. <i>arXiv preprint arXiv:2404.16994</i> .	1525
		1526
		1527
		1528
		1529
	Ruyi Xu, Yuan Yao, Zonghao Guo, Junbo Cui, Zanlin Ni, Chunjiang Ge, Tat-Seng Chua, Zhiyuan Liu, Maosong Sun, and Gao Huang. 2024b. <b>Llava-uhd: an lmm perceiving any aspect ratio and high-resolution images</b> . <i>ArXiv</i> , abs/2403.11703.	1530
		1531
		1532
		1533
		1534
	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. 2024. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> .	1535
		1536
		1537
		1538
	Linli Yao, Lei Li, Shuhuai Ren, Lean Wang, Yuanxin Liu, Xu Sun, and Lu Hou. 2024. Deco: Decoupling token compression from semantic abstraction in multimodal large language models. <i>arXiv preprint arXiv:2405.20985</i> .	1539
		1540
		1541
		1542
		1543
	Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Guohai Xu, Chenliang Li, Junfeng Tian,	1544
		1545

1546	Qi Qian, Ji Zhang, et al. 2023a. Ureader: Universal ocr-free visually-situated language understanding with multimodal large language model. <i>arXiv preprint arXiv:2310.05126</i> .	1601
1547		1602
1548		1603
1549		1604
1550	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:2307.02469</i> .	1605
1551	1606	
1552	1607	
1553	1608	
1554	1609	
1555	1610	
1556	1611	
1557	1612	
1558	1613	
1559	1614	
1560	1615	
1561	1616	
1562	1617	
1563	1618	
1564	1619	
1565	1620	
1566	1621	
1567	1622	
1568	1623	
1569	1624	
1570	1625	
1571	1626	
1572	1627	
1573	1628	
1574	1629	
1575	1630	
1576	1631	
1577	1632	
1578	1633	
1579	1634	
1580	1635	
1581	1636	
1582	1637	
1583	1638	
1584	1639	
1585	1640	
1586	1641	
1587	1642	
1588	1643	
1589	1644	
1590	1645	
1591	1646	
1592	1647	
1593	1648	
1594	1649	
1595	1650	
1596	1651	
1597	1652	
1598	1653	
1599	1654	
1600	1655	
	1656	



1657 Xin Zhang, Dong Zhang, Shimin Li, Yaqian Zhou, and  
1658 Xipeng Qiu. 2023e. Speeche tokenizer: Unified speech  
1659 tokenizer for speech large language models. *arXiv  
1660 preprint arXiv:2308.16692*.

1661 Han Zhao, Min Zhang, Wei Zhao, Pengxiang Ding,  
1662 Siteng Huang, and Donglin Wang. 2024a. Cobra: Ex-  
1663 tending mamba to multi-modal large language model  
1664 for efficient inference.

1665 Han Zhao, Min Zhang, Wei Zhao, Pengxiang Ding,  
1666 Siteng Huang, and Donglin Wang. 2024b. Co-  
1667 bra: Extending mamba to multi-modal large lan-  
1668 guage model for efficient inference. *arXiv preprint  
1669 arXiv:2403.14520*.

1670 Duo Zheng, Shijia Huang, Lin Zhao, Yiwu Zhong, and  
1671 Liwei Wang. 2024a. Towards learning a generalist  
1672 model for embodied navigation. In *Proceedings of  
1673 the IEEE/CVF Conference on Computer Vision and  
1674 Pattern Recognition*, pages 13624–13634.

1675 Kaizhi Zheng, Xuehai He, and Xin Eric Wang.  
1676 2024b. Minigpt-5: Interleaved vision-and-language  
1677 generation via generative vokens. *Preprint,  
1678 arXiv:2310.02239*.

1679 Baichuan Zhou, Ying Hu, Xi Weng, Junlong Jia, Jie Luo,  
1680 Xien Liu, Ji Wu, and Lei Huang. 2024a. Tinyllava: A  
1681 framework of small-scale large multimodal models.  
1682 *arXiv preprint arXiv:2402.14289*.

1683 Gengze Zhou, Yicong Hong, and Qi Wu. 2024b.  
1684 Navgpt: Explicit reasoning in vision-and-language  
1685 navigation with large language models. In *Proceed-  
1686 ings of the AAAI Conference on Artificial Intelligence*,  
1687 volume 38, pages 7641–7649.

1688 Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Ci-  
1689 hang Xie, Alan Yuille, and Tao Kong. 2021. ibot:  
1690 Image bert pre-training with online tokenizer. *arXiv  
1691 preprint arXiv:2111.07832*.

1692 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and  
1693 Mohamed Elhoseiny. 2023a. Minigpt-4: Enhancing  
1694 vision-language understanding with advanced large  
1695 language models. *arXiv:2304.10592*.

1696 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and  
1697 Mohamed Elhoseiny. 2023b. Minigpt-4: Enhancing  
1698 vision-language understanding with advanced large  
1699 language models. *arXiv:2304.10592*.

## A More about Visual Encoder

In Section 2.1.2, we illustrate the visual representations obtained via different visual encoders. As shown in Figure 5, the training strategies of these encoders vary. Most visual encoders are pre-trained in supervised or self-supervised manner. For supervised learning, early exploration utilize image categories as supervision signals (Dosovitskiy et al., 2021), while CLIP-like models (Radford et al., 2021; Sun et al., 2023; Zhai et al., 2023) use language supervision to learn generalized representations. Additionally, SAM (Kirillov et al., 2023) leverages segmentation tasks as training objectives. In contrast, self-supervised learning only requires images for training. One line of works employ contrastive self-supervised methods to distinguish representations between different images (He et al., 2020; Caron et al., 2021; Zhou et al., 2021; Oquab et al., 2023). Another line of approaches construct auto-encoders, where models are demanded to reconstruct images from the encoded visual representations, which is often used to support downstream image generation (Van Den Oord et al., 2017; Esser et al., 2021; Ge et al., 2023a; Sun et al., 2024a).

Since most visual encoders are limited to fixed resolutions and capture certain aspect ratios of visual features, existing LMMs proposes to enhance the input visual representations on two aspects: resolution enhancement and feature enhancement.

To support high-resolutional image processing, the direct method is to increase the resolution accepted by the visual encoder, including interpolating position embeddings in vision Transformers (Zhu et al., 2023a; Bai et al., 2023b) and using CNN-based models to enhance the encoding efficiency of high-resolution images while compressing the size of encoded feature maps (Yuan et al., 2024; Ge et al., 2024). Other works propose to crop high-resolution images into multiple sub-images and input them into the low-resolution encoder along with the down-sampled full image (Ye et al., 2023a; Li et al., 2023e; Gao et al., 2024; Xu et al., 2024b; Liu et al., 2024a). Different sub-image partitioning templates also help address issues caused by varying aspect ratios of images.

Regarding feature enhancement, common practices consider to ensemble visual representations encoded by different encoders, such as combining encoders trained with different strategies (Lu et al., 2024a; Zhao et al., 2024a), or integrating high-resolution and low-resolution encoders to-

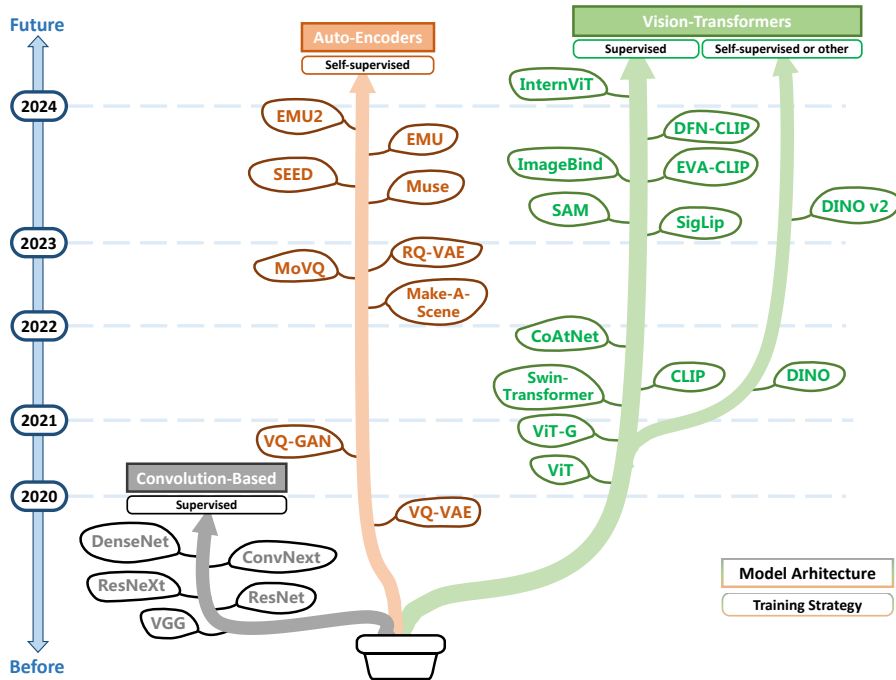


Figure 5: The evolution of commonly adopted visual encoder architectures and training strategies.

gether (Hong et al., 2023a; Li et al., 2024d). Specialized modules have been introduced to better fuse features from different encoders (Li et al., 2024d; Tong et al., 2024; Fan et al., 2024)

## B Multi-Modal Alignment Training

Here, we provide additional details about modality alignment training that could not be fully discussed in the main text due to page limits. In Section 3, we have demonstrated the alignment architectures of contemporary LMMs, as summarized in Table 1 and Table 2. In this section, we will further illustrate the alignment training of LMMs.

The training of current LMMs typically involve multiple stages, with each stage using different data to train specific parameters, gradually learning cross-modal alignment as well as multi-modal understanding and generation capabilities. Most LMMs undergo two main stages: pre-training and instruction fine-tuning. Some models also have additional training stages for specific capabilities.

**Pre-training** The primary goal of pre-training is to align and associate the input representations of various modalities within the multi-modal input space, enabling the backbone to uniformly model and understand inputs across modalities. Figure 6 illustrates the commonly applied settings in the pre-training phase which is described below. At this stage, commonly used data include X-text pairs (“X” means modality X) and multi-modal inter-

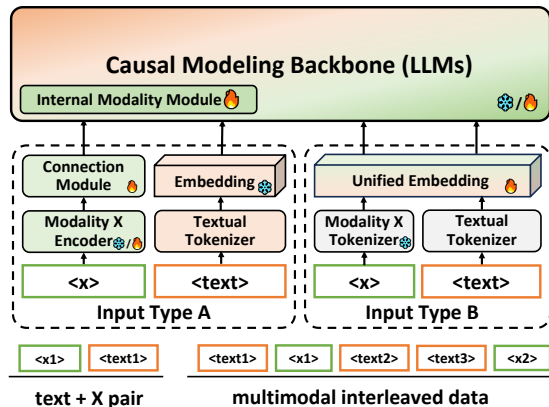


Figure 6: Illustration of common settings during the **pre-training stage**, including data and trainable parameters. “<x>” represents inputs of modalities other than text.

leaved documents. Besides the multi-modal data, text-only data can be adopted to maintain the language modeling capabilities of backbones (Zhang et al., 2023c; Lin et al., 2023; Lu et al., 2024a).

**Instruction Fine-tuning** The instruction fine-tuning stage enables the model to understand and follow instructions to generate appropriate responses, thereby enhancing the interactivity. Figure 7 provides a straightforward illustration for this stage. At this stage, to obtain better generalization under unseen scenarios and tasks, the training data must contain various instructions. Therefore, most LMMs adopt different strategies to construct a mixed dataset based on different re-

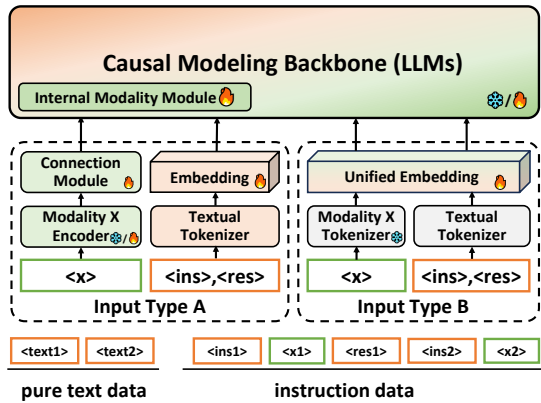


Figure 7: Illustration of common settings during the **instruction fine-tuning stage**, where  $\langle x \rangle$ ,  $\langle \text{ins} \rangle$ , and  $\langle \text{res} \rangle$  denote inputs of modalities other than text, instruction, and response, respectively.

quirements, such as mixing task-oriented data with self-instructed data (Liu et al., 2023; Laurençon et al., 2024b), combining general data with data from specific scenarios (Chen et al., 2023c; Cai et al., 2024), unifying data from various modalities (Wu et al., 2023d; Zhan et al., 2024a; Li et al., 2024b), integrating understanding and generation data (Dong et al., 2024a; Sun et al., 2024a), and blending multi-modal data and text-only data (Lin et al., 2024b; McKinzie et al., 2024).

**Additional Alignment Training** In addition to the regular pre-training and instruction fine-tuning stages, some specialized models require additional training stages to achieve alignment for specific objectives. To enable LLMs to generate multi-modal response, output-level alignment is required. Benefiting from unified multi-modal discrete representation and the pre-trained tokenizer and detokenizer for each modality, models with Type 3 output space can achieve output-level alignment directly through conventional pre-training and instruction fine-tuning (Jin et al., 2023; Ge et al., 2023b; Team, 2024; Zhan et al., 2024a). For models with Type 2 hybrid output space, an additional alignment stage may be required. By rearranging the order of text and other-modality information in “text + X” pairs and interleaved sequences, the text-to-other-modalities generation ability can be learned in the autoregressive setting. A line of approaches keeps modality decoders frozen and train the output mapping modules through gradients passed from the decoder for alignment (Tang et al., 2023c; Dong et al., 2024a; Zheng et al., 2024b). Since most modality decoders are originally conditioned on text for generation, the representations

from the decoders’ corresponding text encoder can be utilized as supervision signal (Wu et al., 2023d; Koh et al., 2024). Another line of methods, represented by Emu series (Sun et al., 2024b,a), propose to construct an autoencoder between modality encoders and decoders. These methods first train LLMs to align the visual input and output spaces, then align the modality decoders to this space.

## C Further Discussions

**How to construct multi-modal representation spaces with discretely or continuously encoded modality signal?** Currently, mainstream LLMs follow the hybrid structure, where modality signals are continuously encoded and integrated into the text space. This method is simple yet effective, leveraging encoders like CLIP (Radford et al., 2021) and CLAP (Elizalde et al., 2023), which are aligned with text through large-scale pre-training, to achieve impressive performance in comprehension tasks. However, this approach introduce additional design costs for corresponding alignment modules for the input and output ends.

Meanwhile, hybrid input spaces cannot directly support multi-modal content generation. This necessitates the design of more complex output layers and decoding strategies for LLMs with multi-modal generation capabilities, leading to a significant gap between the input and output spaces.

On the other hand, the unified discrete space structure is more straightforward, supporting both comprehension and generation tasks through a unified approach (e.g., next-token prediction). However, they are currently limited by the absence of strong discrete encoders across various modalities, akin to CLIP, resulting in slightly weaker performance on comprehension tasks compared to hybrid models. Ovis (Lu et al., 2024b), however, has shown that by carefully designing and expanding the visual vocabulary, discrete models can also perform well on comprehension tasks. Additionally, due to the competitive relationship between modalities, improving training stability is also a challenge that needs to be addressed for unified discrete representation models.

In conclusion, both approaches have their strengths and weaknesses, with significant room for optimization. At the same time, we believe that the current training strategies of discrete and continuous encoders are not mutually exclusive, the development and approaches of both methods can



Model	Input Space		Output Space		Architecture				Max Res.	Date
	Modality	Type	Modality	Type	Backbone	Modality Encoder	Connection	Internal Module		
Flamingo (2022)	Text, Vision	A	Text	1	Chinchilla	NFNet	Perceiver	Cross-Attention	480	2022/04
BLIP-2 (2023b)	Text, Vision	A	Text	1	Flan-T5 / OPT	CLIP ViT-L/14 / Eva-CLIP ViT-G/14	Q-Former	-	224	2023/01
LLaMA-adapter (2023d)	Text, Vision	A	Text	1	LLaMA	CLIP-ViT-L/14	MLP	Adaption Prompt	224	2023/03
MiniGPT-4 (2023b)	Text, Vision	A	Text	1	Vicuna	Eva-CLIP ViT-G/14	Q-Former	-	224	2023/04
LLaVA (2024b)	Text, Vision	A	Text	1	Vicuna	CLIP ViT-L/14	Linear	-	224	2023/04
mPLUG-Owl (2023b)	Text, Vision	A	Text	1	LLaMA	CLIP ViT-L/14	Abstractor	-	224	2023/04
LLaMA-adapter v2 (2023)	Text, Vision	A	Text	1	LLaMA	CLIP-ViT-L/14	MLP	Adaption Prompt	224	2023/04
InstructBLIP (2023)	Text, Vision	A	Text	1	Flan-T5 / Vicuna	Eva-CLIP ViT-G/14	Q-Former	-	224	2023/05
Otter (2023a)	Text, Vision	A	Text	1	LLaMA	CLIP ViT-L/14	Perceiver	Cross-Attention	224	2023/05
LAVIN (2023a)	Text, Vision	A	Text	1	LLaMA	CLIP ViT-L/14	MLP	MM-Adapter	224	2023/05
MultimodalGPT (2023)	Text, Vision	A	Text	1	LLaMA	CLIP ViT-L/14	Perceiver	Cross-Attention	224	2023/05
Shikra (2023c)	Text, Vision	A	Text	1	Vicuna	CLIP ViT-L/14	Linear	-	224	2023/06
VideoChatGPT (2023)	Text, Vision	A	Text	1	Vicuna	CLIP ViT-L/14	Linear	-	224	2023/06
Valley (2023b)	Text, Vision	A	Text	1	Stable-Vicuna	CLIP ViT-L/14	Temporal Module + Linear	-	224	2023/06
Lynx (2023)	Text, Vision	A	Text	1	Vicuna	EVA-1B	Resampler	Adapter	420	2023/07
Qwen-VL (2023b)	Text, Vision	A	Text	1	Qwen	OpenCLIP ViT-bigG	Cross-Attention	-	448	2023/08
BLIVA (2023)	Text, Vision	A	Text	1	Flan-T5 / Vicuna	Eva-CLIP ViT-G/14	Q-Former + MLP	-	224	2023/08
IDEFICS (2023)	Text, Vision	A	Text	1	LLaMA	OpenCLIP ViT-H/14	Perceiver	Cross-Attention	224	2023/08
OpenFlamingo (2023)	Text, Vision	A	Text	1	LLaMA, MPT	CLIP ViT-L/14	Perceiver	Cross-Attention	224	2023/08
InterLM-XC (2023c)	Text, Vision	A	Text	1	InternLM	Eva-CLIP ViT-G/14	Perceiver	-	224	2023/09
LLaVA-1.5 (2023)	Text, Vision	A	Text	1	Vicuna 1.5	CLIP ViT-L/14	MLP	-	336	2023/10
MiniGPT-v2 (2023b)	Text, Vision	A	Text	1	LLaMA-2	EVA	Linear	-	448	2023/10
Fuyu-SB (2023)	Text, Vision	A	Text	1	Persimmon	-	Linear	-	unlimited	2023/10
UReader (2023a)	Text, Vision	A	Text	1	LLaMA	CLIP ViT-L/14	Abstractor	-	224*20	2023/10
CogVLM (2023a)	Text, Vision	A	Text	1	Vicuna 1.5	EVA2-CLIP-E	MLP	Visual Expert	490	2023/11
Monkey (2023e)	Text, Vision	A	Text	1	Qwen	OpenCLIP ViT-bigG	Cross-Attention	-	896	2023/11
ShareGPT4V (2023d)	Text, Vision	A	Text	1	Vicuna-1.5	CLIP ViT-L/14	MLP	-	336	2023/11
mPLUG-Owl2 (2023c)	Text, Vision	A	Text	1	LLaMA-2	CLIP ViT-L/14	Abstractor	Modality-Adaptive Module	448	2023/11
Sphinx (2023)	Text, Vision	A	Text	1	LLaMA-2	CLIP ViT-L/14 + CLIP ConvNeXt-XXL + DINOv2 ViT-G/14	Linear + Q-Former	-	672	2023/11
InternVL (2023f)	Text, Vision	A	Text	1	Vicuna	InternViT	QLLaMA / MLP	-	336	2023/12
MobileVLM (2023a)	Text, Vision	A	Text	1	MobileLLaMA	CLIP ViT-L/14	LDP (conv-based)	-	336	2023/12
VILA (2024b)	Text, Vision	A	Text	1	LLaMA-2	CLIP ViT-L	Linear	-	336	2023/12
Osprey (2024)	Text, Vision	A	Text	1	Vicuna	CLIP ConvNeXt-L	MLP	-	512	2023/12
Honeybee (2024)	Text, Vision	A	Text	1	Vicuna-1.5	CLIP ViT-L/14	C-Abstractor / D-Abstractor	-	336	2023/12
Omni-SMoLA (2024a)	Text, Vision	A	Text	1	UL2	Siglip ViT-G/14	Linear	LoRA MoE	1064	2023/12
LLaVA-Next (2024a)	Text, Vision	A	Text	1	Vicuna / Mistral / Hermes-2-Yi	CLIP ViT-L/14	MLP	-	672	2024/01
InterLM-XC2 (2024b)	Text, Vision	A	Text	1	InternLM-2	CLIP ViT-L/14	MLP	Partial LoRA	490	2024/01
Mousi (2024)	Text, Vision	A	Text	1	Vicuna-1.5	CLIP ViT-L/14 + MAE + LayoutLMv3 + ConvNeXt + SAM + DINOv2 ViT-G	Poly-Expert Fusion	-	1024	2024/01
LLaVA-MoE (2024b)	Text, Vision	A	Text	1	Vicuna 1.5	CLIP ViT-L/14	MLP	LoRA MoE	336	2024/01
MoE-LLaVA (2024a)	Text, Vision	A	Text	1	StableL / Qwen / Phi-2	CLIP ViT-L/14	MLP	FFN MoE	336	2024/01
MobileVLM v2 (2024a)	Text, Vision	A	Text	1	MobileLLaMA	CLIP ViT-L/14	LDP v2	-	336	2024/02
Bunny (2024)	Text, Vision	A	Text	1	Phi-1.5 / LLaMA-3 StableLM-2 / Phi-2	SigLIP, EVA-CLIP	MLP	-	1152	2024/02
TinyLLaVA (2024a)	Text, Vision	A	Text	1	TinyLLaMA / Phi-2 / StableLM-2	SigLIP-L, CLIP ViT-L	MLP	-	336/384	2024/02
Sphinx-X (2024)	Text, Vision	A	Text	1	TinyLLaMA / InternLM2 / LLaMA2 / Mixtral	CLIP ConvNeXt-XXL + DINOv2 ViT-G/14	Linear	-	672	2024/02
Mini-Gemini (2024d)	Text, Vision	A	Text	1	Gemma / Vicuna / Mixtral / Hermes-2-Yi	CLIP ViT-L + ConvNext-L	Cross-Attention + MLP	-	1536	2024/03
Deepseek-VL (2024a)	Text, Vision	A	Text	1	Deepseek LLM	SigLIP-L, SAM-B	MLP	-	1024	2024/03
LLaVA-UHD (2024b)	Text, Vision	A	Text	1	Vicuna	CLIP ViT-L/14	Perceiver	-	336*6	2024/03
Yi-VL (2024)	Text, Vision	A	Text	1	Yi	CLIP ViT-H/14	MLP	-	448	2024/03
MM1 (2024)	Text, Vision	A	Text	1	in-house LLM	CLIP ViT-H*	C-Abstractor	-	1792	2024/03
VL Mamba (2024)	Text, Vision	A	Text	1	Mamba LLM	CLIP-ViT-L / SigLIP-SO400M	VSS + MLP	-	384	2024/03
Cobra (2024b)	Text, Vision	A	Text	1	Mamba-Zephyr	DINOv2 + SigLIP	MLP	-	384	2024/03
InternVL 1.5 (2024c)	Text, Vision	A	Text	1	InternLM2	InternViT-6B	MLP	-	448*40	2024/04
Phi-3-Vision (2024)	Text, Vision	A	Text	1	Phi-3	CLIP ViT-L/14	MLP	-	336*16	2024/04
PLLaVA (2024a)	Text, Vision	A	Text	1	Vicuna / Mistral / Hermes-2-Yi	CLIP ViT-L/14	MLP + Adaptive Pooling	-	336	2024/04
TextHawk (2024a)	Text, Vision	A	Text	1	InternLM-1	SigLIP-SO400M/14	Resampler + MLP	-	unlimited	2024/04
Imp (2024)	Text, Vision	A	Text	1	Phi-2	SigLIP	MLP	-	384	2024/05
IDEFICS2 (2024b)	Text, Vision	A	Text	1	Mistral-v0.1	SigLIP-SO400M/14	Perceiver + MLP	-	384*4	2024/05
ConvLLaVA (2024)	Text, Vision	A	Text	1	Vicuna-LLaMA3	CLIP-ConvNeXt-L*	MLP	-	1536	2024/05
Ovis (2024b)	Text, Vision	B	Text	1	LLaMA3 / Qwen1.5	CLIP ViT-L + Visual Embedding	-	-	336	2024/05
Deco (2024)	Text, Vision	A	Text	1	Vicuna-1.5	CLIP ViT-L/14	MLP + Adaptive Pooling	-	336	2024/05
CuMo (2024c)	Text, Vision	A	Text	1	Mistral / Mixtral	CLIP ViT-L/14	MLP	FFN + MLP MoE	336	2024/05
Cambrian-1 (2024)	Text, Vision	A	Text	1	Vicuna-1.5 / LLaMA-3 / Hermes-2-Yi	CLIP ViT-L/14 + DINOv2 ViT-L/14 + SigLIP ViT-SO400M + OpenCLIP ConvNeXt-XXL	Spatial Vision Aggregator	-	1024	2024/06
GLM-4v (2024)	Text, Vision	A	Text	1	GLM4	EVA-CLIP-E	Conv + SwiGLU	-	1120	2024/06
InterLM-XC2.5 (2024b)	Text, Vision	A	Text	1	InternLM-2	CLIP ViT-L/14	MLP	Partial LoRA	560*24	2024/07
IDEFICS3 (2024a)	Text, Vision	A	Text	1	LLaMA 3.1	SigLIP-SO400M/14	Perceiver + MLP	-	1820	2024/08
mPLUG-Owl3 (2024)	Text, Vision	A	Text	1	Qwen2	SigLIP-SO400M/14	Linear	Hyper Attention	384*6	2024/08
CogVLM2 (2024)	Text, Vision	A	Text	1	LLaMA3	EVA-CLIP-E	Conv + SwiGLU	Visual Expert	1344	2024/08
CogVLM2-video (2024)	Text, Vision	A	Text	1	LLaMA3	EVA-CLIP-E	Conv + SwiGLU	-	224	2024/08
LLaVA-OV (2024a)	Text, Vision	A	Text	1	Qwen-2	SigLIP-SO400M/14	MLP	-	384*36	2024/09
Qwen2-VL (2024)	Text, Vision	A	Text	1	Qwen-2	ViT-675M	MLP	-	unlimited	2024/09

Table 1: Summary of various frameworks of LVLMs that focus on understanding tasks with only text output (Output Type 1). If there are multiple components in a column, ‘+’ represents a combination while ‘/’ indicates an either-or choice. Max Res. represents the maximum resolution, the “X\*Y” pattern indicates methods based on sub-image tiling, X is the base resolution while Y is the maximum number of tiles.

Model	Input Space		Output Space		Architecture						Date
	Modality	Type	Modality	Type	Backbone	Modality Encoder	Connection	Internal Module	Mapping	Modality Decoder	
Any-Modality LLMs											
PandaGPT (2023)	T, V, A...	A	T	1	Vicuna	ImageBind	Linear	-	-	-	2023/05
ImageBind-LLM (2023)	T, V, A, 3D	A	T	1	Chinese-LLaMA	ImageBind + Point-Bind	Bind Network	Adaption Prompt	-	-	2023/09
Next-GPT (2023d)	T, V, A	A	T, V, A	2	Vicuna	ImageBind	Linear	-	Transformer	SD + AudioLDM + Zeriscope	2023/09
Codi-2 (2023c)	T, V, A	A	T, V, A	2	LLaMA-2	ImageBind	MLP	-	MLP	SD + AudioLDM2 + zeroscope v2	2023/11
UnifiedIO2 (2023a)	T, V, A	A	T, V, A	3	UnifiedIO2	OpenCLIP ViT-B + AST	Linear + Perceiver	-	-	VQ-GAN + ViT-VQGAN	2023/12
AnyGPT (2024a)	T, V, A	B	T, V, A	3	LLaMA-2	SEED + Encodec + SpeechTokenizer	-	-	-	SEED + Encodec + SpeechTokenizer	2024/02
Uni-MoE (2024e)	T, V, A	A	T	1	LLaMA	CLIP ViT-L/14 + Whisper-small + BEATs	MLP + Q-former	Modality Aware FFN MoE	-	-	2024/05
Large Audio-Language Models											
SpeechGPT (2023a)	T, A	B	T, A	3	LLaMA	HuBERT	-	-	-	Unit Vocoder	2023/05
Speech-LLaMA (2023b)	T, A	A	T	1	LLaMA	CTC compressor	Transformer	-	-	-	2023/07
SALMONN(2023a)	T, A	A	T	1	Vicuna	Whisper-Large-v2 + BEATs	Window-level Q-Former	-	-	-	2023/10
Qwen-Audio(2023b)	T, A	A	T	1	Qwen	Whisper-Large-v2	-	-	-	-	2023/11
SpeechGPT-Gen (2024a)	T, A	B	T, A	3	LLaMA-2	SpeechTokenizer	-	-	Flow Matching	SpeechTokenizer	2024/01
SLAM-ASR (2024)	T, A	A	T	1	LLaMA-2	HuBERT	MLP + DownSample	-	-	-	2024/02
WavLLM (2024b)	T, A	A	T	1	LLaMA-2	Whisper-Large-v2 + WavLM-Base	Adapter + Linear	-	-	-	2024/04
SpeechVerse (2024)	T, A	A	T	1	Flan-T5-XL	WavLM-Large / Best-RQ	Convolution	-	-	-	2024/05
Qwen2-Audio (2024b)	T, A	A	T	1	Qwen	Whisper-Large-v3	-	-	-	-	2024/07
LLaMA-Omni (2024)	T, A	A	T, A	2	LLaMA-3.1	Whisper-Large-v3	MLP + DownSample	-	Transformer	Unit Vocoder	2024/09
Large Vision-Language Models for Multi-Modal Generation											
GILL (2024)	T, V	A	T, V	2	OPT	CLIP ViT-L	Linear	-	Transformer	SD	2023/05
Emu (2024b)	T, V	A	T, V	2	LLaMA	EVA-02-CLIP-1B	Transformer	-	Linear	SD	2023/07
LaVIT (2023)	T, V	A	T, V	3	LLaMA	Eva-CLIP ViT-G/14 + LaVIT Tokenizer	Linear	-	-	LaVIT De-Tokenizer	2023/09
CM3Leon (2023)	T, V	B	T, V	3	CM3Leon	Make-A-Scene	-	-	-	Make-A-Scene	2023/09
DreamLLM (2024a)	T, V	A	T, V	2	Vicuna	CLIP ViT-L/14	Linear	-	Linear	SD	2023/09
Kosmos-G (2024)	T, V	A	T, V	2	MAGNETO	CLIP ViT-L/14	Resampler	-	AlignerNet	SD	2023/10
SEED-LLaMA (2023b)	T, V	B	T, V	3	Vicuna / LLaMA-2	SEED Tokenizer	-	-	-	SEED De-Tokenizer	2023/10
MiniGPT-5 (2024b)	T, V	A	T, V	2	Vicuna	Eva-CLIP ViT-G/14	Q-Former	-	Transformer	SD	2023/10
Emu-2 (2024a)	T, V	A	T, V	2	LLaMA	EVA-02-CLIP-E-plus	Linear	-	Linear	SDXL	2023/12
Chameleon (2024)	T, V	B	T, V	3	Chameleon	Make-A-Scene	-	-	-	Make-A-Scene	2024/05
MoMA (2024c)	T, V	B	T, V	3	Chameleon	Make-A-Scene	-	Modality Aware FFN MoE	-	Make-A-Scene	2024/07
Vila-U (2024b)	T, V	B	T, V	3	LLaMA-2	SigLIP + RQ-VAE	-	-	-	RQ-VAE	2024/09

Table 2: Supplement to Table 1. In the ‘‘Modality’’ column, T, V, A and 3D are abbreviations for text, vision, audio, and 3D point cloud, respectively.

learn from each other. The research community eagerly anticipates an effective modality encoding method that unifies understanding and generation.

Furthermore, there is a noticeable granularity gap between textual and modal representations, whether the modality signals are encoded continuously or discretely. Text tokens carry explicit semantics, while individual modality tokens might only contain limited information. A single text token may correspond to multiple tokens in an image, leading to excessively long token sequences for modality signals in current LLMs. In the future, can we build modality representations that carry semantics at specific levels?

**How to design model architectures to align the constructed multi-modal space?** The architectures should be designed based on the input and output space. Most LLMs are built on a backbone, usually initialized from a pre-trained LLM to gain better text understanding capabilities and initial rep-

resentations. For hybrid spaces, additional design is required for input and output alignment modules. Although the LLM backbone can perform unified multi-modal modeling through training, relatively complex internal alignment modules can be introduced to model complex cross-modal interactions.

As introduced in Section 3, there is a variety of designs for each module, with different structures having trade-offs across various dimensions. No structure consistently performs better across different scenarios and requirements. Finding ways to quickly validate the effectiveness of an optimization direction is essential. Luckily, there have already been relevant explorations to provide some general conclusions (Laurençon et al., 2024b; McKinzie et al., 2024), offering heuristic approaches to narrow down the model design space.