

From Specific-MLLMs to Omni-MLLMs: A Survey on MLLMs Aligned with Multi-modalities

Anonymous ACL submission

Abstract

To tackle complex tasks in real-world scenarios, more researchers are focusing on Omni-MLLMs, which aim to achieve omni-modal understanding and generation. Beyond the constraints of any specific non-linguistic modality, Omni-MLLMs map various non-linguistic modalities into the embedding space of LLMs and enable the interaction and understanding of arbitrary combinations of modalities within a single model. In this paper, we systematically investigate relevant research and provide a comprehensive survey of Omni-MLLMs. Specifically, we first explain the four core components of Omni-MLLMs for unified multi-modal modeling with a meticulous taxonomy that offers novel perspectives. Then, we introduce the effective integration achieved through two-stage training and discuss the corresponding datasets as well as evaluation. Furthermore, we summarize the main challenges of current Omni-MLLMs and outline future directions. We hope this paper serves as an introduction for beginners and promotes the advancement of related research. Resources will be made public.

1 Introduction

The remarkable performance of continuously evolving Multi-modal Large Language Models (MLLMs) has pointed to a possible direction for achieving general artificial intelligence (Bubeck et al., 2023; OpenAI, 2023b). MLLMs extend Large Language Models (LLMs) by integrating them with pre-trained models tailored to specific modalities, such as Vision-MLLMs (Liu et al., 2023c; Wang et al., 2024b; Sun et al., 2024c), Audio-MLLMs (Zhang et al., 2023a; Chu et al., 2023), and 3D-MLLMs (Xu et al., 2024b). However, these modality-specific MLLMs (Specific-MLLMs) are insufficient to tackle complex tasks in real-world scenarios that simultaneously involve multiple modalities. Therefore, efforts are being made to expand the range of modalities for

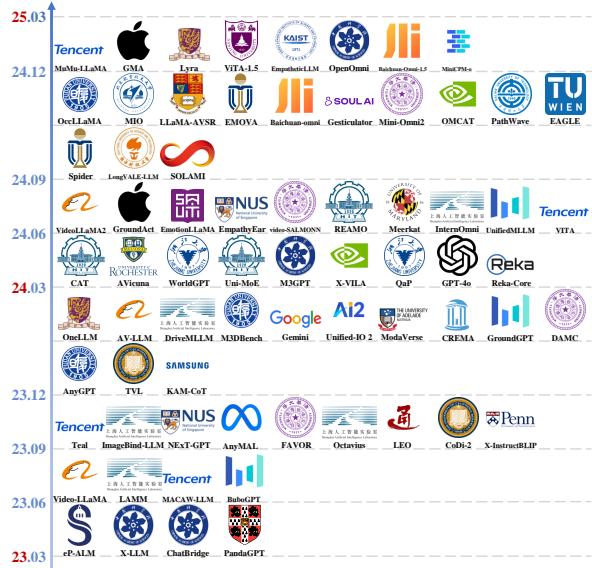


Figure 1: The timeline of representative Omni-MLLMs.

understanding and generation, giving rise to the omni-modality MLLMs (Omni-MLLMs).

By integrating multiple pre-trained models of more non-linguistic modalities (Radford et al., 2021, 2023; Xue et al., 2024b; Rombach et al., 2022; Liu et al., 2023b), Omni-MLLMs expand the modalities for understanding and even generation based on Specific-MLLMs. Omni-MLLMs leverage the emergent capabilities of LLMs to treat various non-linguistic modalities as different *foreign languages*, enabling the interaction and understanding of information across different modalities within a unified space (Chen et al., 2023a; Panagopoulou et al., 2024). Compared to Specific-MLLMs, Omni-MLLMs can perform multiple unimodal¹ understanding and generation tasks, as well as cross-modal tasks across two or more non-linguistic modalities, allowing a single model to

¹“Uni” and “Cross” refer to the number of non-linguistic modalities involved in the interaction, in contrast to “multi-modal reasoning,” traditionally reserved for vision-language tasks (Panagopoulou et al., 2024).

handle arbitrary combinations of modalities.

A review of the development of Omni-MLLMs reveals that it has been continuously expanding in three directions. On the one hand, the types of modalities processed by Omni-MLLM have been continuously increasing, from X-LLMs that handle vision and audio to X-InstructBLIP (Panagopoulou et al., 2024) which adds 3D modality capabilities, PandaGPT (Su et al., 2023) that incorporates IMU modality, and finally One-LLM (Han et al., 2024a), which processes eight different modalities simultaneously. On the other hand, the ability to interact across modalities of Omni-MLLMs has also expanded, from the joint 3D-Image and Audio-Image cross-modal reasoning capability in ImageBind-LLM (Han et al., 2023) to the cross-modal generation capability of CoDi-2 that leverages interleaved audio and image contexts to generate both audio and images (Tang et al., 2024b). The Omni-MLLM is thus trending towards an “Any-to-Any” model. Besides, the application scenarios of Omni-MLLMs have been broadened, encompassing real-time multimodal speech interaction like Mini-Omni2 and Lyra (Xie and Wu, 2024; Zhong et al., 2024), world simulation like WordGPT (Ge et al., 2024b), multi-sensor autonomous driving like DriveMLM (Wang et al., 2023b), etc. In addition to the open-source models, there are also some closed-source Omni-MLLMs such as GPT-4o (OpenAI), Gemini (Reid et al., 2024), and Reka (Ormazabal et al., 2024). The timeline of Omni-MLLMs is shown in Figure 1. Despite the emergence of numerous Omni-MLLMs, there is still a lack of systematic evaluation and analysis.

To fill the gap, we propose this work to conduct a comprehensive and detailed analysis of Omni-MLLMs. We first review the architecture of Omni-MLLMs in four parts (§2). Next, we summarize how Omni-MLLMs expand across multiple modalities through the two-stage training process (§3); then present the training data construction and performance evaluation (§4). Furthermore, we highlight some key challenges and future directions (§5). Finally, we provide a brief summary (§6) and discuss related surveys in the Appendix A.

Our contributions can be summarized as follows: (1) **Comprehensive Survey**: This is the first comprehensive survey dedicated for Omni-MLLMs; (2) **Meticulous taxonomy**: We introduce a meticulous taxonomy (shown in Figure 2); (3) **Challenges and Future**: We outline the challenges of Omni-MLLMs and shed light on future research.

2 Omni-MLLM Architecture

As the extension of Specific-MLLMs, Omni-MLLMs inherit the architecture of *encoding*, *alignment*, *interaction*, and *generation* and broaden the types of non-linguistic modalities involved. This section introduces the implementation methods and functions of the four components in Omni-MLLM: Multi-modalities Encoding (§2.1), Multi-modalities Alignment (§2.2), Multi-modalities Interaction (§2.3), and Multi-modalities Generation (§2.4). More details about the architecture of Omni-MLLMs are shown in Appendix B.

2.1 Multi-modalities Encoding

Based on the encoding feature spaces of multiple modalities, we categorize the Omni-MLLM encoding methods into three types: 1) continuous encoding, 2) discrete encoding, and 3) hybrid encoding.

2.1.1 Continuous Encoding

Continuous encoding refers to encoding the modality into the continuous feature space. Omni-MLLMs that adopt continuous encoding, such as X-LLM (Chen et al., 2023a) and ChatBridge (Zhao et al., 2023b), often integrate multiple pre-trained uni-modality encoders. These modality-specific encoders encode different modalities \mathbf{X} into distinct feature spaces \mathbb{R}_x as \mathbf{F}_x , formulated as:

$$\mathbf{F}_x = \text{SpecificEncoder}(\mathbf{X}), \quad \mathbf{F}_x \in \mathbb{R}_x \quad (1)$$

where SpecificEncoder refers to different modality-specific encoders used in Omni-MLLMs, such as InternVit (Chen et al., 2023g) for encoding visual modality, Whisper (Radford et al., 2023) for encoding auditory modality, ULIP-2 (Xue et al., 2024b) for encoding 3D modality, IMU2CLIP (Moon et al., 2022) for encoding IMU modality, etc.

Besides using heterogeneous encoders for continuous encoding, some Omni-MLLMs (Han et al., 2024a, 2023; Su et al., 2023; Fu et al., 2024c) employ pre-aligned encoders for multiple modalities, encoding different modalities \mathbf{X} into the same feature space \mathbb{R}_{uni} , as shown in Equation 2.

$$\mathbf{F}_x = \text{PreAlignEncoder}(\mathbf{X}), \quad \mathbf{F}_x \in \mathbb{R}_{uni} \quad (2)$$

where PreAlignEncoder refer to encoders that uniformly encode multiple modalities, such as LanguageBind (Zhu et al., 2024a) which uses text as a bridge to align different modalities, and ImageBind (Girdhar et al., 2023) which uses images as a bridge to align different modalities.

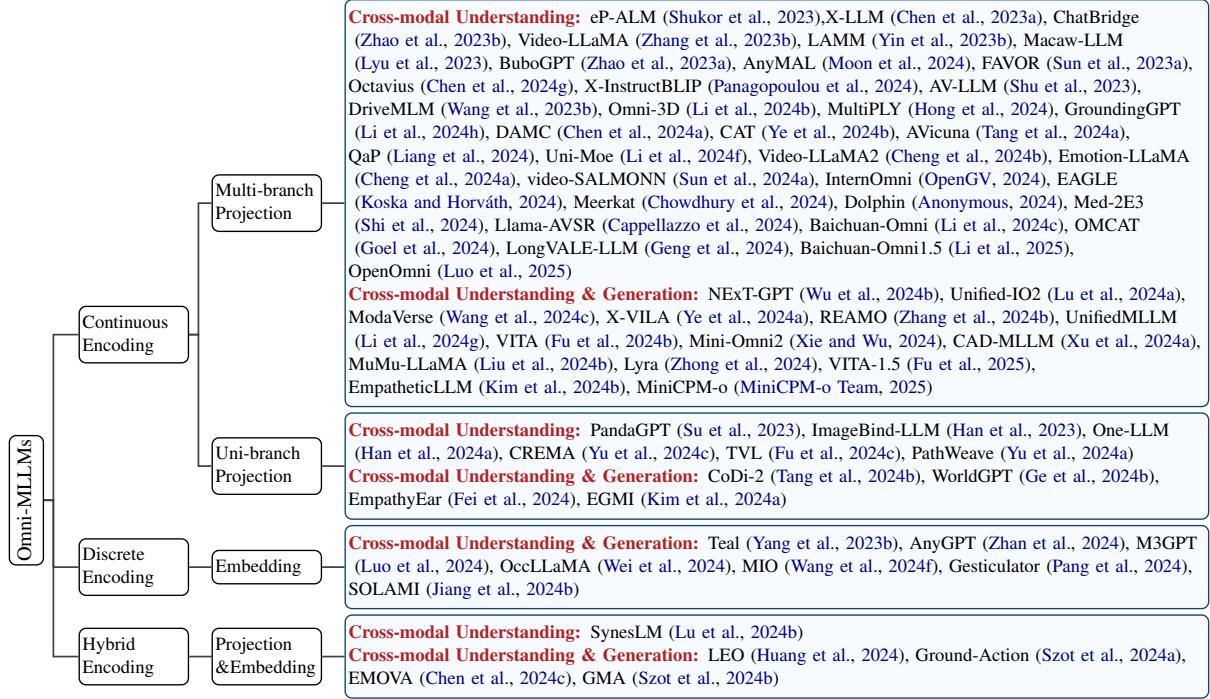


Figure 2: Taxonomy for Omni-MLLMs based on their encoding and alignment methods.

2.1.2 Discrete Encoding

To better facilitate the seamless integration and generation of new non-linguistic modalities, some Omni-MLLMs, such as AnyGPT (Zhan et al., 2024) and Teal (Yang et al., 2023b), adopt a discrete encoding approach. This method encodes different raw modalities \mathbf{X} into the same discrete token space \mathbb{V}_{uni} as \mathbf{T}_x , formulated as follows:

$$\mathbf{T}_x = \text{SpecificTokenizer}(\mathbf{X}), \quad \mathbf{T}_x \in \mathbb{V}_{uni} \quad (3)$$

where SpecificTokenizer refers to different modality-specific tokenizers used in Omni-MLLMs, including the SEED tokenizer (Ge et al., 2024a) based on Vector Quantized Tokenization (VQ), the SpeechTokenizer (Zhang et al., 2023c) based on Residual Vector Quantized Tokenization (RVQ), the AudioTokenizer of Teal (Yang et al., 2023b) based on k-means clustering, and so on.

2.1.3 Hybrid Encoding

Although discrete encoding facilitates the unified processing of different non-linguistic modalities and text compared to continuous encoding, discrete modality tokens often struggle to capture the detailed information inherent in raw continuous modalities (Chen et al., 2024c; Xie and Wu, 2024). Therefore, some Omni-MLLMs combine both encoding approaches instead of a fully discretized manner, choosing different encoding methods for

different modalities. For instance, EMOVA (Chen et al., 2024c) uses the discrete S2U tokenizer to encode auditory modalities while employing the continuous encoder InternVit for visual modalities to retain more vision semantic information. Similarly, GroundAction (Szot et al., 2024a) encodes visual modalities using the CLIP Vit and action modalities with its trained action tokenizer.

2.2 Multi-modalities Alignment

Omni-MLLMs align the encoded features of various non-linguistic modalities with the embedding space of LLMs. The multi-modality alignment can be categorized into two approaches: 1) projection alignment and 2) embedding alignment.

2.2.1 Projection Alignment

The continuous encoding Omni-MLLMs insert adapters, referred to as *projectors*, between the encoders and the LLMs. These projectors map the continuously encoded modality features \mathbf{F}_x into the text embedding space as \mathbf{F}_p . As discussed in Section 2.1.1, \mathbf{F}_x may either reside in distinct feature spaces \mathbb{R}_x or share the same feature space \mathbb{R}_{uni} . For the former, multiple projectors are typically employed to align the \mathbf{F}_x of each modality into \mathbb{R}_t as \mathbf{F}_p independently, addressing dimensional mismatch and feature misalignment across modalities (Ye et al., 2024a; Lyu et al., 2023; Moon et al.,

214), formulated as follows:

$$215 \quad \mathbf{F}_p = \text{SpecificProjector}(\mathbf{F}_x), \mathbf{F}_p \in \mathbb{R}_t \quad (4)$$

216 where SpecificProjector refers to the modality-
217 specific projector corresponding to different modal-
218 ities, called *multi-branch projection*.

219 For the latter case, besides the multi-branch ap-
220 proach, Omni-MLLMs like PandaGPT (Su et al.,
221 2023) and WorldGPT (Ge et al., 2024b) adopt
222 a shared projector to achieve unified alignment
223 across modalities to reduce the parameters of mul-
224 tiple projectors, as shown in Equation 5.

$$225 \quad \mathbf{F}_p = \text{UnifiedProjector}(\mathbf{F}_x), \mathbf{F}_p \in \mathbb{R}_t \quad (5)$$

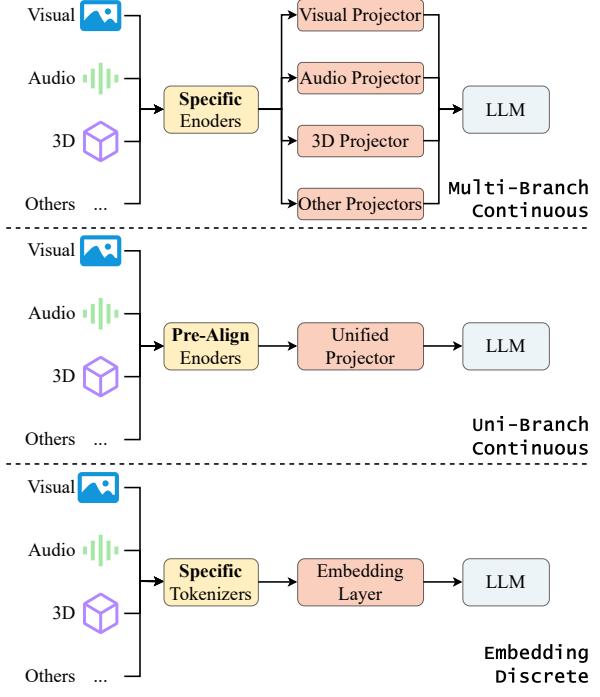
226 where UnifiedProjector refers to the unified pro-
227 jector used to align multiple modalities, a design
228 known as the *uni-branch projection*. A comparison
229 of the two approaches is illustrated in Figure 3.

230 In terms of the **specific implementation** of the
231 projector, the most straightforward approach is
232 to use a multi-layer perceptron (MLP) or a sin-
233 gle linear layer (Wu et al., 2024b; Cheng et al.,
234 2024b; OpenGV, 2024). Alternatively, attention
235 mechanisms can be employed to compress the
236 encoded information of non-linguistic modalities.
237 This includes cross-attention-based methods like
238 Q-Former (Panagopoulou et al., 2024; Chen et al.,
239 2023a) and Perceiver (Zhao et al., 2023b; Liang
240 et al., 2024), as well as self-attention-based meth-
241 ods such as UPM in OneLLM (Han et al., 2024a).
242 Additionally, BaiChuan-Omni (Li et al., 2024c) and
243 EMOVA (Chen et al., 2024c) incorporate CNNs to
244 compress the projected features, thereby achieving
245 locality preservation (Cha et al., 2024).

246 It is also worth noting that in multi-branch Omni-
247 MLLMs, different branches may utilize distinct
248 implementations to better accommodate the unique
249 characteristics of each modality (Li et al., 2024h).
250 For example, Uni-MoE (Li et al., 2024f) uses a
251 linear projection for the visual modality and a Q-
252 Former for the auditory modality. Meanwhile, uni-
253 branch Omni-MLLMs, when using an attention-
254 based projector, typically design multiple modality-
255 specific learnable vectors to extract key information
256 from various non-linguistic modalities (Yu et al.,
257 2024a; Han et al., 2024a; Yu et al., 2024c).

258 2.2.2 Embedding Alignment

259 As for discrete encoding Omni-MLLMs, the fea-
260 tures of non-linguistic modalities are represented
261 as quantized codes, which reside in the same dis-
262 crete space \mathbb{V}_{uni} as text tokens. Therefore, new



263 Figure 3: The three combinations of encoding and align-
264 ment in Omni-MLLM are based on different encoding
265 spaces and alignment structures.

266 modality-specific discrete tokens \mathbf{T}_x are embedded
267 into the continuous feature space \mathbb{R}_t by modifying
268 the vocabulary of LLMs and the corresponding
269 embeddings layer, as shown in Equation 6.

$$270 \quad \mathbf{F}_p = \text{Embedding}(\mathbf{T}_x), \mathbf{F}_p \in \mathbb{R}_t \quad (6)$$

271 where Embedding refers to the unified embedding
272 layer corresponding to different modalities, which
273 is typically achieved by adding discrete codebooks
274 from various modalities to the vocabulary and ex-
275 panding the embedding layer of LLMs (Zhan et al.,
276 2024; Yang et al., 2023b; Wei et al., 2024). For
277 instance, AnyGPT extends the vocabulary of the
278 LLaMA-2 by incorporating 17,408 codes across
279 three modalities—image, speech, and music (Zhan
280 et al., 2024). Besides, some works like Ground-
281 Action (Szot et al., 2024a) and LEO (Huang et al.,
282 2024) overwrite infrequently used tokens in the
283 original vocabulary for alignment, as they extend a
284 smaller set of modality-specific discrete tokens.

285 Additionally, for hybrid encoding models, align-
286 ment is achieved by simultaneously employing
287 both the projection method and the embedding
288 method (Chen et al., 2024c; Szot et al., 2024b).

289 2.3 Multi-modalities Interaction

290 Omni-MLLMs utilize transformer-based LLMs to
291 facilitate information interaction between different

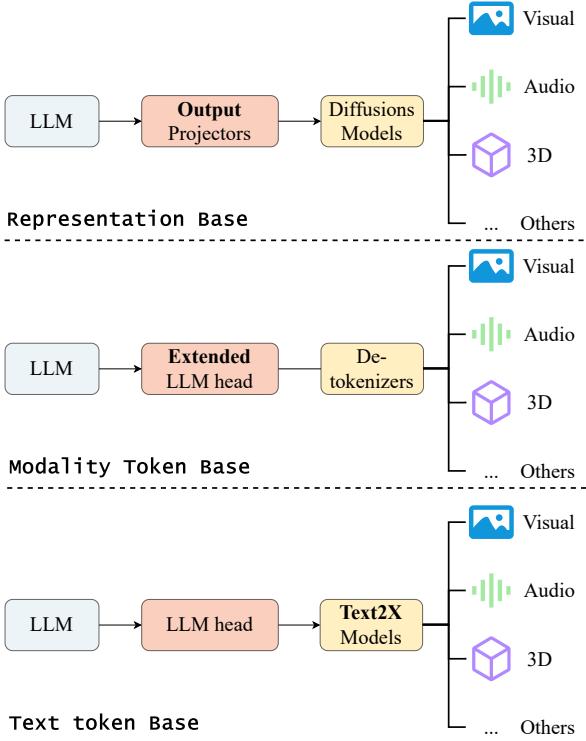


Figure 4: Three generation methods in Omni-MLLMs are implemented based on different output spaces of LLM and the corresponding generative models.

modalities within a unified feature space \mathbb{R}_t . Commonly used LLMs include the LLaMA series (Touvron et al., 2023), the Qwen series (Bai et al., 2023), and others (Cai et al., 2024; Zeng et al., 2024).

For interaction, most Omni-MLLMs (Chen et al., 2023a; Ye et al., 2024a; OpenGV, 2024; Li et al., 2025) concatenate aligned non-linguistic modality features \mathbf{F}_p with textual features \mathbf{F}_t at the input level, enabling interaction in a progressive and layer-by-layer manner. Meanwhile, some works, such as ImageBind-LLM (Han et al., 2023) and TVL-LLaMA (Fu et al., 2024c), insert \mathbf{F}_p into specific layers or all layers of the LLMs to mitigate the loss of modality information (Shukor et al., 2023).

In terms of the number of modalities involved in interactions, compared to Specific-MLLMs that are limited to dual-modal interactions between a single non-linguistic modality and text (Liu et al., 2023c; Xu et al., 2024b), Omni-MLLMs not only support multiple dual-modal interactions but also enable omni-multimodal interactions involving more than two non-linguistic modalities (Zhao et al., 2023b; Wang et al., 2024f). For example, X-InstructBLIP (Panagopoulou et al., 2024) enables dual-modal interactions such as vision-text, audio-text, and 3D-text, as well as omni-modal inter-

actions like vision-audio-text and 3D-vision-text, showcasing the ability of Omni-MLLMs to handle arbitrary combinations of modalities.

2.4 Multi-modalities Generation

Omni-MLLMs can output text while also generating non-linguistic modalities by integrating different generation models. As shown in Figure 4, we categorize multi-modalities generation into three types: text-based generation, representation-based generation, and modality-token-based generation.

Text-based This approach directly utilizes the discrete text output from the LLM to invoke Text-to-X generation models (Liu et al., 2023b; Luo et al., 2023c; Brooks et al., 2023) based on the content of the text. For example, VITA (Fu et al., 2024b) employs TTS tools (RVC-Boss) to convert the output text into corresponding speech, while ModelVerse (Wang et al., 2024c) and UnifiedMLLM (Li et al., 2024g) use the text to specify the generation model and utilize the corresponding descriptions to generate different modalities.

Modality-Token-based Works like MiniOmni-2 (Xie and Wu, 2024) and AnyGPT (Zhan et al., 2024) extend the corresponding LLM head with codebooks from different modality tokenizers to generate modality-specific discrete tokens. These tokens are then decoded using the corresponding de-tokenizers (Esser et al., 2021; Yu et al., 2024b; Zeghidour et al., 2022; Dhariwal et al., 2020) to produce various modalities.

Representation-based To alleviate the potential noise introduced by discrete tokens, works like X-VILA (Ye et al., 2024a) and NextGPT (Wu et al., 2024b) incorporate modality-specific signal tokens into the vocabulary. They then use transformers or MLPs to map the signal token representations into the ones that are understandable to the multimodal decoders, typically off-the-shelf latent-conditioned diffusion models (Rombach et al., 2022; Tang et al., 2023; Xue et al., 2024a; Blattmann et al., 2023), enabling effective generation capabilities.

3 Omni-MLLM Training

To achieve alignment across different vector spaces and improve instruction-following ability under arbitrary modality settings, Omni-MLLMs extend the standard two-stage training pipeline of Specific-MLLMs: *multi-modalities alignment pre-training* and *multi-modalities instruction fine-tuning*.

363 3.1 Multi-modalities Alignment Pre-training

364 Multi-modalities alignment pre-training involves
365 *input alignment* training between the feature spaces
366 of different modalities and the embedding space
367 of LLMs on the encoding side, as well as *output*
368 *alignment* training between the embedding space
369 and the input spaces of various modality decoders
370 on the decoding side. Input alignment and output
371 alignment can be carried out separately (Wu et al.,
372 2024b) or simultaneously (Ye et al., 2024a).

373 3.1.1 Input Alignment

374 Input alignment mainly uses X-Text paired datasets
375 of different modalities and minimizes the text genera-
376 tion loss of the corresponding description text to
377 optimize. In this phase, continuous encoding Omni-
378 MLLMs normally update parameters of projectors,
379 while discrete encoding Omni-MLLMs adjust the
380 parameters of the embedding layer.

381 In terms of *training order* of different modalities
382 alignment, multi-branch Omni-MLLM performs
383 separate alignment training for each modality-
384 specific projector, directly aligning each non-
385 linguistic modality with text and using text as
386 a bridge to align different non-linguistic modalities
387 (Zhao et al., 2023b; Panagopoulou et al., 2024).
388 The uni-branch Omni-MLLM, on the other hand,
389 uses the unified projector for different modalities,
390 which may lead to interference in the alignment
391 performance between different modalities. Thus,
392 Han et al. (2024a) employ a progressive alignment
393 strategy to align multiple modalities in a specific
394 order. In contrast, discrete encoding Omni-MLLMs,
395 like AnyGPT (Zhan et al., 2024) and M3GPT (Luo
396 et al., 2024), mix the alignment data from different
397 modalities and perform alignment simultaneously.

398 Besides, in addition to directly leveraging X-
399 Text paired datasets from different modalities for
400 direct alignment, PandaGPT (Su et al., 2023),
401 ImageBind-LLM (Han et al., 2023), and VideoL-
402 LaMA (Zhang et al., 2023b) utilize the pre-aligned
403 modality feature space \mathbb{V}_{uni} to achieve indirect
404 alignment between other non-linguistic modalities
405 and text by training solely on Image-Text data.

406 3.1.2 Output Alignment

407 The training of output alignment typically utilizes
408 the same X-text paired dataset as input alignment
409 and adheres to the identical training sequence.
410 Meanwhile, the training objectives for output align-
411 ment vary depending on the multi-modality gener-
412 ation methods in Section 2.4. Token-based gener-

413 ative Omni-MLLMs optimize the extended LLM
414 head by minimizing the text generation loss asso-
415 ciated with modality-specific discrete tokens (Lu
416 et al., 2024a; Wei et al., 2024). Representation-
417 based generative Omni-MLLMs generally optimize
418 their output projectors by minimizing the compos-
419 itive loss comprising three components (Xie and Wu,
420 2024; Yang et al., 2023b): 1) the text generation
421 loss of signal tokens; 2) the L2 distance between
422 the output representation and the condition vec-
423 tor of the corresponding decoder, i.e. MSE loss;
424 and 3) the conditional latent denoising loss (Rom-
425 bach et al., 2022). For text-based generative Omni-
426 MLLMs, as there is no additional output structure,
427 the output alignment training is generally not re-
428 quired (Wang et al., 2024c; Li et al., 2024g).

429 3.2 Multi-modalities Instruction Fine-tuning

430 The instruction fine-tuning phase aims to enhance
431 generalization capability under arbitrary modalities
432 of Omni-MLLMs (Panagopoulou et al., 2024;
433 Wu et al., 2024b; Ye et al., 2024a). Instruction
434 fine-tuning primarily utilizes instruction-following
435 datasets and computes the text generation loss for
436 the corresponding responses to optimize. For mod-
437 els with generation capabilities, the loss mentioned
438 in section 3.1.2 may also be incorporated. Dur-
439 ing this phase, Omni-MLLMs further perform full-
440 scale tuning of the LLM parameters (Han et al.,
441 2024a; Cheng et al., 2024b) or use PEFT tech-
442 niques (Han et al., 2024b), such as LoRA (Hu et al.,
443 2022), for partial tuning (Wang et al., 2024e).

444 Compared to Specific-MLLMs, Omni-MLLMs
445 not only leverage multiple uni-modal instruction
446 data of different modalities for training but also use
447 cross-modal instruction data to enhance their cross-
448 modal ability (Ye et al., 2024a; Zhan et al., 2024;
449 Li et al., 2024c). In addition to directly mixing dif-
450 ferent instruction data for training (Panagopoulou
451 et al., 2024; Ye et al., 2024a), some works like Uni-
452 Moe (Li et al., 2024f) and Lyra (Zhong et al., 2024)
453 adopt a multi-step fine-tuning approach, introduc-
454 ing different uni-modal and cross-modal instruc-
455 tion data in a specific order for training to gradually
456 enhance their uni-modal and cross-modal ability.

457 4 Data Construction and Evaluation

458 This section summarizes the construction of modal-
459 ity alignment data and instruction data used in the
460 Omni-MLLM training process (§4.1), as well as the
461 evaluation across four different capabilities (§4.2).

4.1 Training Data

Alignment Data Omni-MLLMs leverage caption datasets from various modalities to construct X-Text paired data for alignment pre-training, such as the WebVid (Bain et al., 2021) for visual modality and the AudioCaps (Kim et al., 2019) for auditory modality. However, for data-scarce modalities like depth maps and thermal maps, large-scale text-paired data is lacking (Zhu et al., 2024a; Girdhar et al., 2023). To address this, synthetic methods that use DPT models (Ranftl et al., 2021; Bhat et al., 2023; Xu et al., 2023) or image translation models (Lee et al., 2023) to convert image-text pairs into other modality text pairs are widely employed (Han et al., 2024a; Chen et al., 2024c; Zhu et al., 2024a). Moreover, interleaved datasets (Zhu et al., 2023a) are used for alignment pre-training in some works (Tang et al., 2024b) to enhance the contextual understanding capability.

Instruction Data Omni-MLLMs not only leverage uni-modal instruction datasets from Specific-MLLMs, but also construct cross-modal instruction data through diverse methods as follows.

(1) **Template-based Construction:** Most works (Sun et al., 2023a; Zhao et al., 2023b; Zhang et al., 2024b) utilize cross-modal downstream datasets (Sanabria et al., 2018; Chen et al., 2020b) combined with predefined templates to construct cross-modal instructions; (2) **GPT Generation:** Following the paradigm of LLaVA (Liu et al., 2023c), some Omni-MLLMs (Lyu et al., 2023; Zhao et al., 2023b) leverage the labels from the annotated dataset (Lin et al., 2014; Bain et al., 2021) or use pre-trained models like SAM (Chen et al., 2023c) and GRIT (Wu et al., 2024a) to extract meta-information (e.g., captions and object categories) of different modalities. Then they employ powerful LLMs (OpenAI, 2023b,a) to generate cross-modal instructions based on the obtained meta-information; (3) **T2X Generation:** Li et al. (2024f) use TTS tools to convert the Image-Text2Text uni-modal instructions from LLaVA-v1.5 (Liu et al., 2024a) into Image-Speech-Text2Text cross-modal instructions. AnyGPT (Zhan et al., 2024) and NextGPT (Wu et al., 2024b) leverage Text2X models such as DALL-E-3 (Shi et al., 2020) and MusicGen (Copet et al., 2023) to convert the GPT-generated pure text instructions into Xs2Xs cross-modal instructions. Details about training data are shown in Appendix C.1

4.2 Benchmark

We provide a brief overview of the benchmarks used to evaluate Omni-MLLMs. The statistics of the benchmarks are shown in Appendix C.2.

Uni-modal Understanding Uni-modal understanding assesses the ability of Omni-MLLMs to comprehend and reason on different non-linguistic modalities, including downstream X-Text2Text datasets such as X-Caption (Plummer et al., 2015; Xu et al., 2016), X-QA (Goyal et al., 2017; Xu et al., 2017), and X-Classification (Deitke et al., 2023), as well as comprehensive multi-task benchmarks (Liu et al., 2024d; Fu et al., 2024a, 2023).

Uni-modal Generation Uni-modal generation aims to evaluate the ability of Omni-MLLMs to generate a single non-linguistic modality, including the Text2X generation task (Kim et al., 2019; Ruiz et al., 2023) and the Text-X2Text editing task (Veaux et al., 2017; Perazzi et al., 2016).

Cross-modal Understanding Cross-modal understanding evaluates the ability of Omni-MLLMs to jointly comprehend and reason across multiple non-linguistic modalities like Image-Speech-Text2Text (Li et al., 2024f; OpenGV, 2024), Video-Audio-Text2Text (Li et al., 2022a,b), and Image-3D-Text2Text (Panagopoulou et al., 2024).

Cross-modal Generation Cross-modal generation further evaluates the ability of Omni-MLLMs to generate non-linguistic modalities in conjunction with other non-linguistic modality inputs. For example, the Xs-Text2X benchmark proposed by X-VILA (Ye et al., 2024a) includes tasks such as Image-Text2Audio and Image-Audio-Text2Video.

5 Challenges and Future Directions

Despite Omni-MLLMs having showcased remarkable performance on numerous tasks, there are still some challenges that necessitate further research.

5.1 Expansion of modalities

Most Omni-MLLMs can only process 2-3 types of non-linguistic modalities, and they still face several challenges when expanding more modalities.

Training efficiency The common method that introduces new modalities through additional alignment pre-training and instruction fine-tuning can lead to significant training cost. Leveraging prior knowledge from Specific-MLLMs (Panagopoulou

et al., 2024; Chen et al., 2024a) or using pre-aligned encoders for indirect alignment (Han et al., 2023; Su et al., 2023) can help reduce training overhead but may impact cross-modal performance.

Catastrophic forgetting Expanding new modalities may adjust the shared parameters, potentially causing catastrophic forgetting of previously trained modalities knowledge (Yu et al., 2024a). This issue can be partially mitigated by mixing trained modality data (Han et al., 2024a; Li et al., 2024f) or fine-tuning only the modality-specific parameters (Yu et al., 2024a,c), but both approaches make the training process more complex.

Low-resource modalities Although the data synthesis method in Section 4.1 can help alleviate the lack of text-paired data and instruction data for low-resource modalities (Han et al., 2024a; ?), the absence of real modality may lead to biases in understanding of that modality.

5.2 Cross-modal capabilities

The Omni-MLLMs have achieved promising performance in cross-modal understanding and generation tasks, but there are still some challenges.

Long Context When the input contains multiple sequence modalities (video, speech...), the length of the multi-modalities token sequence may exceed the context window of LLMs and lead to memory overflow. While methods such as token compressing (Yu et al., 2024c; Li et al., 2024c) or token sampling (Zhan et al., 2024; Zhong et al., 2024) can reduce the number of input tokens, they also result in a decline in cross-modal performance.

Modality Bias Due to the imbalance in training data volume and the performance disparity among different modality encoders, Omni-MLLMs may tend to pay attention to the dominant modality while neglecting information from other modalities during cross-modal inference. Balancing the data volume across modalities or enhancing the corresponding modality-specific modules could potentially help mitigate this issue (Leng et al., 2024).

Temporal Alignment When dealing with different modalities that have temporal dependencies, retaining their temporal alignment information is crucial for subsequent cross-modal understanding. Some attempts have been made to preserve the temporal alignment information between audio and video, such as interleaved modality-specific tokens

of video and audio (Tang et al., 2024a) and inserting the time-related special tokens into the multi-modalities tokens (Goel et al., 2024).

Data and Benchmark Although Omni-MLLMs employ various methods in Section 4.1 to generate cross-modal instruction data, there is still significant room for improvement and expansion, including enhancing the diversity of instructions, incorporating longer contextual dialogues, and exploring more diverse modality interaction paradigms. Similarly, cross-modal benchmarks such as OmniBench (Li et al., 2024e) and OmniR (Chen et al., 2024d) still fall short in terms of task richness and instruction diversity when compared to uni-modal benchmarks like MMMU (Yue et al., 2024) and MME (Fu et al., 2023). And the variety of modalities they cover is also relatively limited.

5.3 Application scenarios

The emergence of Omni-MLLM brings new opportunities and possibilities for various applications. **(1) Real-time Multi-modalities Interaction:** Fu et al. (2025) and Xie and Wu (2024) achieve robust capabilities in both vision and speech understanding, enabling efficient speech-to-speech interactions with vision in real-time. **(2) Comprehensive Planning:** Wang et al. (2023b) and Szot et al. (2024a) leverage the complementarity across multiple modalities to achieve better path planning and action planning capabilities than planning with vision information only. **(3) World Simulator:** Ge et al. (2024b) not only understands and generates different modalities but also predicts state transitions for any combination of modalities.

6 Conclusion

In this paper, we provide a comprehensive survey report on Omni-MLLM, offering a comprehensive review of the field. Specifically, we break down Omni-MLLM into four key components and categorize them based on modal encoding and alignment methods. Subsequently, we provide a detailed summary of the training process of Omni-MLLM and the related resources used. We also summarize the current challenges and the future development directions. This paper is the first systematic survey dedicated to Omni-MLLMs. We hope this survey will facilitate further research in this area.

652 Limitations

653 This study provides the first comprehensive survey
654 of Omni-MLLMs. Related work, architecture
655 statistics, more details of training and evaluation,
656 as well as other training recipes, can be found in
657 Appendix A,B,C.

658 We have made our best effort, but there may
659 still be some limitations. On one hand, due to
660 page limitations, we can only provide a concise
661 overview of the core contributions of mainstream
662 Omni-MLLMs, rather than exhaustive technical
663 details. On the other hand, our review primarily cov-
664 ers research from *ACL, NeurIPS, ICLR, ICML,
665 COLING, CVPR, IJCAI, ECCV, and arXiv, and
666 there is a chance that we may have missed some
667 important work published in other venues. We will
668 stay updated with ongoing discussions in the re-
669 search community and plan to revise our work in
670 the future to include overlooked contributions.

671 References

672 Panos Achlioptas, Ahmed Abdelreheem, Fei Xia, Mo-
673 hamed Elhoseiny, and Leonidas J. Guibas. 2020.
674 *Referit3d: Neural listeners for fine-grained 3d ob-
675 ject identification in real-world scenes*. In *Computer
676 Vision - ECCV 2020 - 16th European Conference,
677 Glasgow, UK, August 23-28, 2020, Proceedings, Part
678 I*, volume 12346 of *Lecture Notes in Computer Sci-
679 ence*, pages 422–440. Springer.

680 Andrea Agostinelli, Timo I. Denk, Zalán Borsos,
681 Jesse H. Engel, Mauro Verzetti, Antoine Caillon,
682 Qingqing Huang, Aren Jansen, Adam Roberts, Marco
683 Tagliasacchi, Matthew Sharifi, Neil Zeghidour, and
684 Christian Havnø Frank. 2023. *Musiclm: Generating
685 music from text*. *CoRR*, abs/2301.11325.

686 Harsh Agrawal, Peter Anderson, Karan Desai, Yufei
687 Wang, Xinlei Chen, Rishabh Jain, Mark Johnson,
688 Dhruv Batra, Devi Parikh, and Stefan Lee. 2019.
689 *nocaps: novel object captioning at scale*. In *2019
690 IEEE/CVF International Conference on Computer
691 Vision, ICCV 2019, Seoul, Korea (South), October 27
692 - November 2, 2019*, pages 8947–8956. IEEE.

693 Huda AlAmri, Vincent Cartillier, Raphael Gontijo
694 Lopes, Abhishek Das, Jue Wang, Irfan Essa, Dhruv
695 Batra, Devi Parikh, Anoop Cherian, Tim K. Marks,
696 and Chiori Hori. 2018. *Audio visual scene-
697 aware dialog (AVSD) challenge at DSTC7*. *CoRR*,
698 abs/1806.00525.

699 Emily J Allen, Ghislain St-Yves, Yihan Wu, Jesse L
700 Breedlove, Jacob S Prince, Logan T Dowdle,
701 Matthias Nau, Brad Caron, Franco Pestilli, Ian
702 Charest, et al. 2022. A massive 7t fmri dataset to
703 bridge cognitive neuroscience and artificial intelli-
704 gence. *Nature neuroscience*, 25(1):116–126.

705 Anonymous. 2024. *Aligned better, listen better for
706 audio-visual large language models*. In *Submitted to
707 The Thirteenth International Conference on Learning
708 Representations*. Under review.

709 Rosana Ardila, Megan Branson, Kelly Davis, Michael
710 Kohler, Josh Meyer, Michael Henretty, Reuben
711 Morais, Lindsay Saunders, Francis M. Tyers, and
712 Gregor Weber. 2020. *Common voice: A massively-
713 multilingual speech corpus*. In *Proceedings of The
714 12th Language Resources and Evaluation Confer-
715 ence, LREC 2020, Marseille, France, May 11-16,
716 2020*, pages 4218–4222. European Language Re-
717 sources Association.

718 Anurag Arnab, Mostafa Dehghani, Georg Heigold,
719 Chen Sun, Mario Lucic, and Cordelia Schmid.
720 2021. *Vivit: A video vision transformer*. In *2021
721 IEEE/CVF International Conference on Computer
722 Vision, ICCV 2021, Montreal, QC, Canada, October
723 10-17, 2021*, pages 6816–6826. IEEE.

724 Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Mo-
725 toaki Kawanabe. 2022. *Scanqa: 3d question answer-
726 ing for spatial scene understanding*. In *IEEE/CVF
727 Conference on Computer Vision and Pattern Recog-
728 nition, CVPR 2022, New Orleans, LA, USA, June
729 18-24, 2022*, pages 19107–19117. IEEE.

730 Fan Bai, Yuxin Du, Tiejun Huang, Max Qinghu Meng,
731 and Bo Zhao. 2024a. *M3D: advancing 3d medical
732 image analysis with multi-modal large language mod-
els*. *CoRR*, abs/2404.00578.

733 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
734 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
735 Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,
736 Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu,
737 Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren,
738 Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong
739 Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang
740 Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian
741 Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi
742 Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang,
743 Yichang Zhang, Zhenru Zhang, Chang Zhou, Jing-
744 ren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023.
745 *Qwen technical report*. *CoRR*, abs/2309.16609.

746 Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He,
747 Zongbo Han, Zheng Zhang, and Mike Zheng Shou.
748 2024b. *Hallucination of multimodal large language
749 models: A survey*. *CoRR*, abs/2404.18930.

750 Max Bain, Arsha Nagrani, Gü̈l Varol, and Andrew Zis-
751 serman. 2021. *Frozen in time: A joint video and
752 image encoder for end-to-end retrieval*. In *2021
753 IEEE/CVF International Conference on Computer
754 Vision, ICCV 2021, Montreal, QC, Canada, October
755 10-17, 2021*, pages 1708–1718. IEEE.

756 Gedas Bertasius, Heng Wang, and Lorenzo Torresani.
757 2021. *Is space-time attention all you need for video
758 understanding?* In *Proceedings of the 38th Inter-
759 national Conference on Machine Learning, ICML
760 2021, 18-24 July 2021, Virtual Event*, volume 139 of
761

762	<i>Proceedings of Machine Learning Research</i> , pages	813–824. PMLR.	819
763			820
764	Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter		821
765	Wonka, and Matthias Müller. 2023. <i>Zoedepth: Zero-</i>		822
766	<i>shot transfer by combining relative and metric depth</i> .		823
767	<i>CoRR</i> , abs/2302.12288.		824
768	Ali Furkan Biten, Rubèn Tito, Andrés Mafla, Lluís		825
769	Gómez, Marçal Rusiñol, Minesh Mathew, C. V. Jawa-		826
770	har, Ernest Valveny, and Dimosthenis Karatzas. 2019.		827
771	ICDAR 2019 competition on scene text visual ques-		828
772	tion answering. In <i>2019 International Conference on</i>		829
773	<i>Document Analysis and Recognition, ICDAR 2019,</i>		830
774	<i>Sydney, Australia, September 20-25, 2019</i> , pages		
775	1563–1570. IEEE.		
776	Andreas Blattmann, Tim Dockhorn, Sumith Ku-		
777	ial, Daniel Mendelevitch, Maciej Kilian, Dominik		
778	Lorenz, Yam Levi, Zion English, Vikram Voleti,		
779	Adam Letts, Varun Jampani, and Robin Rom-		
780	bach. 2023. <i>Stable video diffusion: Scaling latent</i>		
781	<i>video diffusion models to large datasets.</i> <i>CoRR</i> ,		
782	abs/2311.15127.		
783	Tim Brooks, Aleksander Holynski, and Alexei A. Efros.		
784	2023. <i>Instructpix2pix: Learning to follow image</i>		
785	<i>editing instructions.</i> In <i>IEEE/CVF Conference on</i>		
786	<i>Computer Vision and Pattern Recognition, CVPR</i>		
787	<i>2023, Vancouver, BC, Canada, June 17-24, 2023</i> ,		
788	pages 18392–18402. IEEE.		
789	Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao		
790	Zheng. 2017. <i>AISHELL-1: an open-source man-</i>		
791	<i>darin speech corpus and a speech recognition base-</i>		
792	<i>line.</i> In <i>20th Conference of the Oriental Chapter of</i>		
793	<i>the International Coordinating Committee on Speech</i>		
794	<i>Databases and Speech I/O Systems and Assessment,</i>		
795	<i>O-COCOSDA 2017, Seoul, South Korea, November</i>		
796	1-3, 2017, pages 1–5. IEEE.		
797	Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan,		
798	Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter		
799	Lee, Yin Tat Lee, Yuanzhi Li, Scott M. Lundberg,		
800	Harsha Nori, Hamid Palangi, Marco Túlio Ribeiro,		
801	and Yi Zhang. 2023. <i>Sparks of artificial general</i>		
802	<i>intelligence: Early experiments with GPT-4.</i> <i>CoRR</i> ,		
803	abs/2303.12712.		
804	Davide Caffagni, Federico Cocchi, Luca Barsellotti,		
805	Nicholas Moratelli, Sara Sarto, Lorenzo Baraldi, Mar-		
806	cella Cornia, and Rita Cucchiara. 2024. <i>The revolu-</i>		
807	<i>tion of multimodal large language models: A survey.</i>		
808	In <i>Findings of the Association for Computational Lin-</i>		
809	<i>guistics, ACL 2024, Bangkok, Thailand and virtual</i>		
810	<i>meeting, August 11-16, 2024</i> , pages 13590–13618.		
811	Association for Computational Linguistics.		
812	Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen,		
813	Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi		
814	Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan,		
815	Zhaoye Fei, Yang Gao, Jiaye Ge, Chunya Gu, Yuzhe		
816	Gu, Tao Gui, Aijia Guo, Qipeng Guo, Conghui He,		
817	Yingfan Hu, Ting Huang, Tao Jiang, Penglong Jiao,		
818	Zhenjiang Jin, Zhikai Lei, Jiaxing Li, Jingwen Li,		
	Linyang Li, Shuaibin Li, Wei Li, Yining Li, Hong-		
	wei Liu, Jiangning Liu, Jiawei Hong, Kaiwen Liu,		
	Kuikun Liu, Xiaoran Liu, Chengqi Lv, Haijun Lv,		
	Kai Lv, Li Ma, Runyuan Ma, Zerun Ma, Wenchang		
	Ning, Linke Ouyang, Jiantao Qiu, Yuan Qu, Fukai		
	Shang, Yunfan Shao, Demin Song, Zifan Song, Zhi-		
	hao Sui, Peng Sun, Yu Sun, Huanze Tang, Bin Wang,		
	Guoteng Wang, Jiaqi Wang, Jiayu Wang, Rui Wang,		
	Yudong Wang, Ziyi Wang, Xingjian Wei, Qizhen		
	Weng, Fan Wu, Yingtong Xiong, Xiaomeng Zhao,		
	and et al. 2024. <i>Internlm2 technical report.</i> <i>CoRR</i> ,		
	abs/2403.17297.		
	Umberto Cappellazzo, Minsu Kim, Honglie Chen,		
	Pingchuan Ma, Stavros Petridis, Daniele Falavigna,		
	Alessio Brutti, and Maja Pantic. 2024. <i>Large lan-</i>		
	<i>guage models are strong audio-visual speech recog-</i>		
	<i>nition learners.</i> <i>CoRR</i> , abs/2409.12319.		
	Cerspense. 2023. <i>Zeroscope: Diffusion-based text-to-</i>		
	<i>video synthesis.</i>		
	Junbum Cha, Wooyoung Kang, Jonghwan Mun, and		
	Byungseok Roh. 2024. <i>Honeybee: Locality-</i>		
	<i>enhanced projector for multimodal LLM.</i> In		
	<i>IEEE/CVF Conference on Computer Vision and Pat-</i>		
	<i>tern Recognition, CVPR 2024, Seattle, WA, USA,</i>		
	<i>June 16-22, 2024</i> , pages 13817–13827. IEEE.		
	Soravit Changpinyo, Piyush Sharma, Nan Ding, and		
	Radu Soricut. 2021. <i>Conceptual 12m: Pushing web-</i>		
	<i>scale image-text pre-training to recognize long-tail</i>		
	<i>visual concepts.</i> In <i>IEEE Conference on Computer</i>		
	<i>Vision and Pattern Recognition, CVPR 2021, virtual,</i>		
	<i>June 19-25, 2021</i> , pages 3558–3568. Computer Vi-		
	sion Foundation / IEEE.		
	Chi Chen, Yiyang Du, Zheng Fang, Ziyue Wang, Fuwen		
	Luo, Peng Li, Ming Yan, Ji Zhang, Fei Huang,		
	Maosong Sun, and Yang Liu. 2024a. <i>Model com-</i>		
	<i>position for multimodal large language models.</i> In		
	<i>Proceedings of the 62nd Annual Meeting of the As-</i>		
	<i>sociation for Computational Linguistics (Volume 1:</i>		
	<i>Long Papers), ACL 2024, Bangkok, Thailand, August</i>		
	<i>11-16, 2024</i> , pages 11246–11262. Association for		
	Computational Linguistics.		
	Dave Zhenyu Chen, Angel X. Chang, and Matthias		
	Nießner. 2020a. <i>Scanrefer: 3d object localization in</i>		
	<i>RGB-D scans using natural language.</i> In <i>Computer</i>		
	<i>Vision - ECCV 2020 - 16th European Conference,</i>		
	<i>Glasgow, UK, August 23-28, 2020, Proceedings, Part</i>		
	<i>XX, volume 12365 of Lecture Notes in Computer</i>		
	<i>Science</i> , pages 202–221. Springer.		
	Feilong Chen, Minglun Han, Haozhi Zhao, Qingyang		
	Zhang, Jing Shi, Shuang Xu, and Bo Xu. 2023a. <i>X-</i>		
	<i>LLM: bootstrapping advanced large language models</i>		
	<i>by treating multi-modalities as foreign languages.</i>		
	<i>CoRR</i> , abs/2305.04160.		
	Guoguo Chen, Shuzhou Chai, Guan-Bo Wang, Jiayu		
	Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel		
	Povey, Jan Trmal, Junbo Zhang, Mingjie Jin, San-		
	jeev Khudanpur, Shinji Watanabe, Shuaijiang Zhao,		

876 877 878 879 880 881 882	<p>Wei Zou, Xiangang Li, Xuchen Yao, Yongqing Wang, Zhao You, and Zhiyong Yan. 2021. Gigaspeech: An evolving, multi-domain ASR corpus with 10, 000 hours of transcribed audio. In <i>22nd Annual Conference of the International Speech Communication Association, Interspeech 2021, Brno, Czechia, August 30 - September 3, 2021</i>, pages 3670–3674. ISCA.</p>	933 934
883 884 885 886 887 888	<p>Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, Chao Weng, and Ying Shan. 2023b. Videocrafter1: Open diffusion models for high-quality video generation. <i>CoRR</i>, abs/2310.19512.</p>	939 940 941 942 943 944 945 946
889 890 891 892 893	<p>Hong Chen, Xin Wang, Yuwei Zhou, Bin Huang, Yipeng Zhang, Wei Feng, Houlun Chen, Zeyang Zhang, Siao Tang, and Wenwu Zhu. 2024b. Multi-modal generative AI: multi-modal llm, diffusion and beyond. <i>CoRR</i>, abs/2409.14993.</p>	947 948 949 950 951
894 895 896 897 898 899	<p>Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. 2020b. Vggsound: A large-scale audio-visual dataset. In <i>2020 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2020, Barcelona, Spain, May 4-8, 2020</i>, pages 721–725. IEEE.</p>	952 953 954 955
900 901 902	<p>Jiaqi Chen, Zeyu Yang, and Li Zhang. 2023c. Semantic segment anything. https://github.com/fudan-zvg/Semantic-Segment-Anything.</p>	956 957 958
903 904 905 906 907 908 909 910 911 912	<p>Kai Chen, Yunhao Gou, Runhui Huang, Zhili Liu, Daxin Tan, Jing Xu, Chunwei Wang, Yi Zhu, Yihan Zeng, Kuo Yang, Dingdong Wang, Kun Xiang, Haoyuan Li, Haoli Bai, Jianhua Han, Xiaohui Li, Weike Jin, Nian Xie, Yu Zhang, James T. Kwok, Hengshuang Zhao, Xiaodan Liang, Dit-Yan Yeung, Xiao Chen, Zhenguo Li, Wei Zhang, Qun Liu, Jun Yao, Lanqing Hong, Lu Hou, and Hang Xu. 2024c. EMOVA: empowering language models to see, hear and speak with vivid emotions. <i>CoRR</i>, abs/2409.18042.</p>	961 962
913 914 915 916 917 918	<p>Lichang Chen, Hexiang Hu, Mingda Zhang, Yiwen Chen, Zifeng Wang, Yandong Li, Pranav Shyam, Tianyi Zhou, Heng Huang, Ming-Hsuan Yang, and Boqing Gong. 2024d. Omnixr: Evaluating omni-modality language models on reasoning across modalities. <i>CoRR</i>, abs/2410.12219.</p>	963 964 965 966 967 968
919 920 921 922 923 924 925 926	<p>Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2024e. Sharegpt4v: Improving large multi-modal models with better captions. In <i>Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part XVII</i>, volume 15075 of <i>Lecture Notes in Computer Science</i>, pages 370–387. Springer.</p>	969 970 971 972 973
927 928 929 930 931 932	<p>Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, Wanxiang Che, Xiangzhan Yu, and Furu Wei. 2023d. Beats: Audio pre-training with acoustic tokenizers. In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i>, volume</p>	974 975 976 977 978 979
933 934 935 936 937 938	<p>202 of <i>Proceedings of Machine Learning Research</i>, pages 5178–5193. PMLR.</p> <p>Sihan Chen, Xingjian He, Longteng Guo, Xinxin Zhu, Weining Wang, Jinhui Tang, and Jing Liu. 2023e. VALOR: vision-audio-language omni-perception pre-training model and dataset. <i>CoRR</i>, abs/2304.08345.</p> <p>Sihan Chen, Handong Li, Qunbo Wang, Zijia Zhao, Mingzhen Sun, Xinxin Zhu, and Jing Liu. 2023f. VAST: A vision-audio-subtitle-text omni-modality foundation model and dataset. In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i>.</p> <p>Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao, Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, and Sergey Tulyakov. 2024f. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i>, pages 13320–13331. IEEE.</p> <p>Zeren Chen, Ziqin Wang, Zhen Wang, Huayang Liu, Zhenfei Yin, Si Liu, Lu Sheng, Wanli Ouyang, and Jing Shao. 2024g. Octavius: Mitigating task interference in mllms via lora-moe. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i>. OpenReview.net.</p> <p>Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. 2023g. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. <i>CoRR</i>, abs/2312.14238.</p> <p>Zebang Cheng, Zhi-Qi Cheng, Jun-Yan He, Jingdong Sun, Kai Wang, Yuxiang Lin, Zheng Lian, Xiaojiang Peng, and Alexander Hauptmann. 2024a. Emotion-llama: Multimodal emotion recognition and reasoning with instruction tuning. <i>CoRR</i>, abs/2406.11161.</p> <p>Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, and Lidong Bing. 2024b. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-lmms. <i>CoRR</i>, abs/2406.07476.</p> <p>Chee Kheng Chng and Chee Seng Chan. 2017. Totaltext: A comprehensive dataset for scene text detection and recognition. In <i>14th IAPR International Conference on Document Analysis and Recognition, ICDAR 2017, Kyoto, Japan, November 9-15, 2017</i>, pages 935–942. IEEE.</p> <p>Sanjoy Chowdhury, Sayan Nag, Subhrajyoti Dasgupta, Jun Chen, Mohamed Elhoseiny, Ruohan Gao, and Dinesh Manocha. 2024. MEERKAT: audio-visual large language model for grounding in space and time. In</p>	980 981 982 983 984 985

990	<i>Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXIV, volume 15122 of Lecture Notes in Computer Science</i> , pages 52–70. Springer.	1048
991		1049
992		1050
993		1051
994		1052
995		1053
996		1054
997		1055
998		1056
999		1057
1000		1058
1001		1059
1002		1060
1003		1061
1004		1062
1005		1063
1006		1064
1007		1065
1008		1066
1009		1067
1010		1068
1011		1069
1012		1070
1013		1071
1014		1072
1015		1073
1016		1074
1017		1075
1018		1076
1019		1077
1020		1078
1021		1079
1022		1080
1023		1081
1024		1082
1025		1083
1026		1084
1027		1085
1028		1086
1029		1087
1030		1088
1031		1089
1032		1090
1033		1091
1034		1092
1035		1093
1036		1094
1037		1095
1038		1096
1039		1097
1040		1098
1041		1099
1042		1100
1043		1101
1044	<i>Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXIV, volume 15122 of Lecture Notes in Computer Science</i> , pages 52–70. Springer.	1048
1045		1049
1046		1050
1047		1051
1048	<i>Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXIV, volume 15122 of Lecture Notes in Computer Science</i> , pages 52–70. Springer.	1052
1049		1053
1050		1054
1051		1055
1052		1056
1053		1057
1054		1058
1055		1059
1056		1060
1057		1061
1058		1062
1059		1063
1060		1064
1061		1065
1062		1066
1063		1067
1064		1068
1065		1069
1066		1070
1067		1071
1068		1072
1069		1073
1070		1074
1071		1075
1072		1076
1073		1077
1074		1078
1075		1079
1076		1080
1077		1081
1078		1082
1079		1083
1080		1084
1081		1085
1082		1086
1083		1087
1084		1088
1085		1089
1086		1090
1087		1091
1088		1092
1089		1093
1090		1094
1091		1095
1092		1096
1093		1097
1094		1098
1095		1099
1096		1100
1097		1101

1102	Chaoyou Fu, Yuhang Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, Peixian Chen, Yanwei Li, Shaohui Lin, Sirui Zhao, Ke Li, Tong Xu, Xiawu Zheng, Enhong Chen, Rongrong Ji, and Xing Sun. 2024a. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. <i>CoRR</i> , abs/2405.21075.	CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023, pages 22942–22951. IEEE.	1159
1103			1160
1104			1161
1105			1162
1106			1163
1107			1164
1108			1165
1109			
1110	Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Xiong Wang, Di Yin, Long Ma, Xiawu Zheng, Ran He, Rongrong Ji, Yunsheng Wu, Caifeng Shan, and Xing Sun. 2024b. VITA: towards open-source interactive omni multimodal LLM. <i>CoRR</i> , abs/2408.05211.	Tiantian Geng, Jinrui Zhang, Qingni Wang, Teng Wang, Jinming Duan, and Feng Zheng. 2024. Longvale: Vision-audio-language-event benchmark towards time-aware omni-modal perception of long videos. <i>CoRR</i> , abs/2411.19772.	1166
1111			1167
1112			1168
1113			1169
1114			1170
1115			1171
			1172
1116	Chaoyou Fu, Haojia Lin, Xiong Wang, Yi-Fan Zhang, Yunhang Shen, Xiaoyu Liu, Yangze Li, Zuwei Long, Heting Gao, Ke Li, Long Ma, Xiawu Zheng, Rongrong Ji, Xing Sun, Caifeng Shan, and Ran He. 2025. Vita-1.5: Towards gpt-4o level real-time vision and speech interaction. <i>Preprint</i> , arXiv:2501.01957.	Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. 2023. Imagebind one embedding space to bind them all. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023</i> , pages 15180–15190. IEEE.	1166
1117			1167
1118			1168
1119			1169
1120			1170
1121			1171
			1172
1122	Letian Fu, Gaurav Datta, Huang Huang, William Chung-Ho Panitch, Jaimyn Drake, Joseph Ortiz, Mustafa Mukadam, Mike Lambeta, Roberto Calandra, and Ken Goldberg. 2024c. A touch, vision, and language dataset for multimodal alignment. In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	Arushi Goel, Karan Sapra, Matthieu Le, Rafael Valle, Andrew Tao, and Bryan Catanzaro. 2024. OMCAT: omni context aware transformer. <i>CoRR</i> , abs/2410.12109.	1166
1123			1167
1124			1168
1125			1169
1126			1170
1127			1171
1128			1172
1129			
1130	Yuying Ge, Yixiao Ge, Ziyun Zeng, Xintao Wang, and Ying Shan. 2023. Planting a SEED of vision in large language model. <i>CoRR</i> , abs/2307.08041.	Yuan Gong, Yu-An Chung, and James R. Glass. 2021. AST: audio spectrogram transformer. In <i>22nd Annual Conference of the International Speech Communication Association, Interspeech 2021, Brno, Czechia, August 30 - September 3, 2021</i> , pages 571–575. ISCA.	1166
1131			1167
1132			1168
			1169
1133	Yuying Ge, Sijie Zhao, Ziyun Zeng, Yixiao Ge, Chen Li, Xintao Wang, and Ying Shan. 2024a. Making llama SEE and draw with SEED tokenizer. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	Yuan Gong, Jin Yu, and James R. Glass. 2022. Vocal-sound: A dataset for improving human vocal sounds recognition. In <i>IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022</i> , pages 151–155. IEEE.	1166
1134			1167
1135			1168
1136			1169
1137			1170
1138			1171
			1172
1139	Zhiqi Ge, Hongzhe Huang, Mingze Zhou, Juncheng Li, Guoming Wang, Siliang Tang, and Yueting Zhuang. 2024b. Worldgpt: Empowering LLM as multimodal world model. In <i>Proceedings of the 32nd ACM International Conference on Multimedia, MM 2024, Melbourne, VIC, Australia, 28 October 2024 - 1 November 2024</i> , pages 7346–7355. ACM.	Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In <i>2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017</i> , pages 6325–6334. IEEE Computer Society.	1166
1140			1167
1141			1168
1142			1169
1143			1170
1144			1171
1145			1172
			1173
1146	Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. 2017. Audio set: An ontology and human-labeled dataset for audio events. In <i>2017 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2017, New Orleans, LA, USA, March 5-9, 2017</i> , pages 776–780. IEEE.	Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan, Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray, Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Cartillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano Fragnani, Qichen Fu, Abrham Gebreselasie, Cristina González, James Hillis, Xuhua Huang, Yifei Huang, Wenqi Jia, Leslie Khoo, Jáchym Kolář, Satwik Kottur, Anurag Kumar, Federico Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, Jonathan Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz Puentes, Merey Ramazanova, Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo, Yuchen Wang, Xindi Wu, Takuma Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbeláez, David Crandall, Dima Damen, Giovanni Maria Farinella, Christian Fuegen, Bernard	1166
1147			1167
1148			1168
1149			1169
1150			1170
1151			1171
1152			1172
1153			1173
			1174
1154	Tiantian Geng, Teng Wang, Jinming Duan, Runmin Cong, and Feng Zheng. 2023. Dense-localizing audio-visual events in untrimmed videos: A large-scale benchmark and baseline. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> ,		1166
1155			1167
1156			1168
1157			1169
1158			1170
			1171

1218	Ghanem, Vamsi Krishna Ithapu, C. V. Jawahar, Han-	Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and	1275
1219	byul Joo, Kris Kitani, Haizhou Li, Richard A. New-	Sai Qian Zhang. 2024b. Parameter-efficient fine-	1276
1220	combe, Aude Oliva, Hyun Soo Park, James M. Rehg,	tuning for large models: A comprehensive survey.	1277
1221	Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Anto-	<i>CoRR</i> , abs/2403.14608.	1278
1222	nio Torralba, Lorenzo Torresani, Mingfei Yan, and		
1223	Jitendra Malik. 2022. Ego4d: Around the world in		
1224	3, 000 hours of egocentric video. In <i>IEEE/CVF Con-</i>		
1225	<i>ference on Computer Vision and Pattern Recognition,</i>		
1226	<i>CVPR 2022, New Orleans, LA, USA, June 18-24,</i>		
1227	<i>2022, pages 18973–18990. IEEE.</i>		
1228	Jiaxi Gu, Xiaojun Meng, Guansong Lu, Lu Hou, Niu		
1229	Minzhe, Xiaodan Liang, Lewei Yao, Runhui Huang,		
1230	Wei Zhang, Xin Jiang, Chunjing Xu, and Hang Xu.		
1231	2022. Wukong: A 100 million large-scale chinese		
1232	cross-modal pre-training benchmark. In <i>Advances in</i>		
1233	<i>Neural Information Processing Systems 35: Annual</i>		
1234	<i>Conference on Neural Information Processing Sys-</i>		
1235	<i>tems 2022, NeurIPS 2022, New Orleans, LA, USA,</i>		
1236	<i>November 28 - December 9, 2022.</i>		
1237	Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki		
1238	Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang,		
1239	Zhengdong Zhang, Yonghui Wu, and Ruoming Pang.		
1240	2020. Conformer: Convolution-augmented trans-		
1241	former for speech recognition. In <i>21st Annual Con-</i>		
1242	<i>ference of the International Speech Communication</i>		
1243	<i>Association, Interspeech 2020, Virtual Event, Shang-</i>		
1244	<i>hai, China, October 25-29, 2020, pages 5036–5040.</i>		
1245	<i>ISCA.</i>		
1246	Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio		
1247	César Teodoro Mendes, Allie Del Giorno, Sivakanth		
1248	Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo		
1249	de Rosa, Olli Saarikivi, Adil Salim, Shital Shah,		
1250	Harkirat Singh Behl, Xin Wang, Sébastien Bubeck,		
1251	Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and		
1252	Yuanzhi Li. 2023. Textbooks are all you need. <i>CoRR</i> ,		
1253	abs/2306.11644.		
1254	Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo,		
1255	Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P.		
1256	Bigham. 2018. Vizwiz grand challenge: Answering		
1257	visual questions from blind people. In <i>2018 IEEE</i>		
1258	<i>Conference on Computer Vision and Pattern Recog-</i>		
1259	<i>nition, CVPR 2018, Salt Lake City, UT, USA, June</i>		
1260	<i>18-22, 2018, pages 3608–3617. Computer Vision</i>		
1261	<i>Foundation / IEEE Computer Society.</i>		
1262	Jiaming Han, Kaixiong Gong, Yiyuan Zhang, Jiaqi		
1263	Wang, Kaipeng Zhang, Dahua Lin, Yu Qiao, Peng		
1264	Gao, and Xiangyu Yue. 2024a. Onellm: One frame-		
1265	work to align all modalities with language. In <i>IEEE/CVF Conference on Computer Vision and Pat-</i>		
1266	<i>tern Recognition, CVPR 2024, Seattle, WA, USA,</i>		
1267	<i>June 16-22, 2024, pages 26574–26585. IEEE.</i>		
1268			
1269	Jiaming Han, Renrui Zhang, Wenqi Shao, Peng Gao,		
1270	Peng Xu, Han Xiao, Kaipeng Zhang, Chris Liu, Song		
1271	Wen, Ziyu Guo, Xudong Lu, Shuai Ren, Yafei Wen,		
1272	Xiaoxin Chen, Xiangyu Yue, Hongsheng Li, and		
1273	Yu Qiao. 2023. Imagebind-llm: Multi-modality in-		
1274	struction tuning. <i>CoRR</i> , abs/2309.03905.		
	Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and		
	Sai Qian Zhang. 2024b. Parameter-efficient fine-		
	tuning for large models: A comprehensive survey.		
	<i>CoRR</i> , abs/2403.14608.		
	Yingqing He, Zhao Yang Liu, Jingye Chen, Zeyue Tian,		
	Hongyu Liu, Xiaowei Chi, Runtao Liu, Ruibin Yuan,		
	Yazhou Xing, Wenhui Wang, Jifeng Dai, Yong Zhang,		
	Wei Xue, Qifeng Liu, Yike Guo, and Qifeng Chen.		
	2024. Llms meet multimodal generation and editing:		
	A survey. <i>CoRR</i> , abs/2405.19334.		
	Yining Hong, Zishuo Zheng, Peihao Chen, Yian Wang,		
	Junyan Li, and Chuang Gan. 2024. Multiply: A		
	multisensory object-centric embodied large language		
	model in 3d world. In <i>IEEE/CVF Conference on</i>		
	<i>Computer Vision and Pattern Recognition, CVPR</i>		
	<i>2024, Seattle, WA, USA, June 16-22, 2024, pages</i>		
	<i>26396–26406. IEEE.</i>		
	Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai,		
	Kushal Lakhota, Ruslan Salakhutdinov, and Abdellah-		
	rahman Mohamed. 2021. Hubert: Self-supervised		
	speech representation learning by masked prediction		
	of hidden units. <i>IEEE ACM Trans. Audio Speech</i>		
	<i>Lang. Process.</i> , 29:3451–3460.		
	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan		
	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and		
	Weizhu Chen. 2022. Lora: Low-rank adaptation of		
	large language models. In <i>The Tenth International</i>		
	<i>Conference on Learning Representations, ICLR 2022,</i>		
	<i>Virtual Event, April 25-29, 2022. OpenReview.net.</i>		
	Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun		
	Linghu, Puhaol Li, Yan Wang, Qing Li, Song-Chun		
	Zhu, Baoxiong Jia, and Siyuan Huang. 2024. An		
	embodied generalist agent in 3d world. In <i>Forty-</i>		
	<i>first International Conference on Machine Learning,</i>		
	<i>ICML 2024, Vienna, Austria, July 21-27, 2024. Open-</i>		
	<i>Review.net.</i>		
	Jiaxing Huang and Jingyi Zhang. 2024. A survey		
	on evaluation of multimodal large language models.		
	<i>CoRR</i> , abs/2408.15769.		
	Xiaoshui Huang, Sheng Li, Wentao Qu, Tong He, Yi-		
	fan Zuo, and Wanli Ouyang. 2022. Frozen CLIP		
	model is an efficient point cloud backbone. <i>CoRR</i> ,		
	abs/2212.04098.		
	Drew A. Hudson and Christopher D. Manning. 2019.		
	GQA: A new dataset for real-world visual reason-		
	ing and compositional question answering. In <i>IEEE</i>		
	<i>Conference on Computer Vision and Pattern Recog-</i>		
	<i>nition, CVPR 2019, Long Beach, CA, USA, June 16-20,</i>		
	<i>2019, pages 6700–6709. Computer Vision Founda-</i>		
	<i>tion / IEEE.</i>		
	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-		
	sch, Chris Bamford, Devendra Singh Chaplot, Diego		
	de Las Casas, Florian Bressand, Gianna Lengyel,		
	Guillaume Lample, Lucile Saulnier, Lélio Ren-		
	nard Lavaud, Marie-Anne Lachaux, Pierre Stock,		
	Teven Le Scao, Thibaut Lavril, Thomas Wang, Timo-		
	thée Lacroix, and William El Sayed. 2023. Mistral		
	7b. <i>CoRR</i> , abs/2310.06825.		

1333	Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024a. <i>Mixtral of experts</i> . <i>CoRR</i> , abs/2401.04088.	1391
1334		1392
1335		1393
1336		1394
1337		1395
1338		1396
1339		1397
1340		1398
1341		1399
1342		
1343	Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. 2019. <i>Audiocaps: Generating captions for audios in the wild</i> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)</i> , pages 119–132. Association for Computational Linguistics.	1391
1344		1392
1345		1393
1346		1394
1347		1395
1348		1396
1349		1397
1350	Jianping Jiang, Weiye Xiao, Zhengyu Lin, Huaizhong Zhang, Tianxiang Ren, Yang Gao, Zhiqian Lin, Zhonggang Cai, Lei Yang, and Ziwei Liu. 2024b. <i>SOLAMI: social vision-language-action modeling for immersive interaction with 3d autonomous characters</i> . <i>CoRR</i> , abs/2412.00174.	1400
1351		1401
1352		1402
1353		1403
1354		1404
1355		1405
1356		
1357	Taemin Kim, WOOYEOL BAEK, and Heeseok Oh. 2024a. <i>Efficient generative multimodal integration (EGMI): Enabling audio generation from text-image pairs through alignment with large language models</i> . In <i>Audio Imagination: NeurIPS 2024 Workshop AI-Driven Speech, Music, and Sound Generation</i> .	1400
1358		1401
1359		1402
1360		1403
1361		1404
1362		1405
1363	Yeonju Kim, Se Jin Park, and Yong Man Ro. 2024b. <i>Empathetic response in audio-visual conversations using emotion preference optimization and mamba-compressor</i> . <i>CoRR</i> , abs/2412.17572.	1406
1364		1407
1365		1408
1366		1409
1367	Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D. Plumley. 2020. <i>Panns: Large-scale pretrained audio neural networks for audio pattern recognition</i> . <i>IEEE ACM Trans. Audio Speech Lang. Process.</i> , 28:2880–2894.	1410
1368		1411
1369		1412
1370		1413
1371		1414
1372	Ben Koska and Mojmír Horváth. 2024. <i>Towards multi-modal mastery: A 4.5b parameter truly multi-modal small language model</i> . <i>CoRR</i> , abs/2411.05903.	1415
1373		1416
1374		1417
1375	Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. 2017a. <i>Dense-captioning events in videos</i> . In <i>IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017</i> , pages 706–715. IEEE Computer Society.	1418
1376		1419
1377		1420
1378		1421
1379		1422
1380		1423
1381	Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. 2017b. <i>Visual genome: Connecting language and vision using crowdsourced dense image annotations</i> . <i>Int. J. Comput. Vis.</i> , 123(1):32–73.	1424
1382		1425
1383	Jinxiang Lai, Jie Zhang, Jun Liu, Jian Li, Xiaocheng Lu, and Song Guo. 2024. <i>Spider: Any-to-many multimodal LLM</i> . <i>CoRR</i> , abs/2411.09439.	1426
1384		1427
1385		1428
1386		1429
1387		1430
1388	Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. 2023. <i>OBELICS: an open web-scale filtered dataset of interleaved image-text documents</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1431
1389		1440
1390	Dong-Guw Lee, Myung-Hwan Jeon, Younggun Cho, and Ayoung Kim. 2023. <i>Edge-guided multi-domain rgb-to-tir image translation for training vision tasks</i>	1441
1391		1442
1392		1443
1393		
1394		
1395		
1396		
1397		
1398		
1399		

1447	with challenging labels.	In <i>IEEE International Conference on Robotics and Automation, ICRA 2023, London, UK, May 29 - June 2, 2023</i> , pages 8291–8298. IEEE.	1505
1448			1506
1449			1507
1450			1508
1451	Sicong Leng, Yun Xing, Zesen Cheng, Yang Zhou, Hang Zhang, Xin Li, Deli Zhao, Shijian Lu, Chunyan Miao, and Lidong Bing.	2024. The curse of multi-modalities: Evaluating hallucinations of large multimodal models across language, visual, and audio. <i>CoRR</i> , abs/2410.12787.	1509
1452			1510
1453			1511
1454			1512
1455			1513
1456			1514
1457	Guangyao Li, Yake Wei, Yapeng Tian, Chenliang Xu, Ji-Rong Wen, and Di Hu.	2022a. Learning to answer questions in dynamic audio-visual scenarios. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 19086–19096. IEEE.	1515
1458			1516
1459			1517
1460			1518
1461			1519
1462			1520
1463	Guangyao Li, Yake Wei, Yapeng Tian, Chenliang Xu, Ji-Rong Wen, and Di Hu.	2022b. Learning to answer questions in dynamic audio-visual scenarios. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 19086–19096. IEEE.	1521
1464			1522
1465			1523
1466			1524
1467			1525
1468	Jian Li and Weiheng Lu.	2024. A survey on benchmarks of multimodal large language models. <i>CoRR</i> , abs/2408.08632.	1526
1469			1527
1470			1528
1471			1529
1472	Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi.	2022c. BLIP: bootstrapping language-image pre-training for unified vision-language understanding and generation. In <i>International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA</i> , volume 162 of <i>Proceedings of Machine Learning Research</i> , pages 12888–12900. PMLR.	1530
1473			1531
1474			1532
1475			1533
1476			1534
1477			1535
1478			1536
1479			1537
1480	Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Lou, Limin Wang, and Yu Qiao.	2024a. Mvbench: A comprehensive multi-modal video understanding benchmark. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 22195–22206. IEEE.	1538
1481			1539
1482			1540
1483			1541
1484			1542
1485			1543
1486			1544
1487			1545
1488	Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, Licheng Yu, and Jingjing Liu.	2020. HERO: hierarchical encoder for video+language omni-representation pre-training. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 2046–2065. Association for Computational Linguistics.	1546
1489			1547
1490			1548
1491			1549
1492			
1493			
1494			
1495			
1496	Mingsheng Li, Xin Chen, Chi Zhang, Sijin Chen, Hongyuan Zhu, Fukun Yin, Zhuoyuan Li, Gang Yu, and Tao Chen.	2024b. M3dbench: Towards omni 3d assistant with interleaved multi-modal instructions. In <i>Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LVIII</i> , volume 15116 of <i>Lecture Notes in Computer Science</i> , pages 41–59. Springer.	1550
1497			1551
1498			1552
1499			1553
1500			1554
1501			1555
1502			1556
1503			1557
1504			1558
1505			1559
1506	Yadong Li, Jun Liu, Tao Zhang, Tao Zhang, Song Chen, Tianpeng Li, Zehuan Li, Lijun Liu, Lingfeng Ming, Guosheng Dong, Da Pan, Chong Li, Yuanbo Fang, Dongdong Kuang, Mingrui Wang, Chenglin Zhu, Youwei Zhang, Hongyu Guo, Fengyu Zhang, Yuran Wang, Bowen Ding, Wei Song, Xu Li, Yuqi Huo, Zheng Liang, Shusen Zhang, Xin Wu, Shuai Zhao, Linchu Xiong, Yozhen Wu, Jiahui Ye, Wenhao Lu, Bowen Li, Yan Zhang, Yaqi Zhou, Xin Chen, Lei Su, Hongda Zhang, Fuzhong Chen, Xuezhen Dong, Na Nie, Zhiying Wu, Bin Xiao, Ting Li, Shunya Dang, Ping Zhang, Yijia Sun, Jincheng Wu, Jinjie Yang, Xionghai Lin, Zhi Ma, Kegeng Wu, Jia li, Aiyuan Yang, Hui Liu, Jianqiang Zhang, Xiaoxi Chen, Guangwei Ai, Wentao Zhang, Yicong Chen, Xiaoqin Huang, Kun Li, Wenjing Luo, Yifei Duan, Lingling Zhu, Ran Xiao, Zhe Su, Jiani Pu, Dian Wang, Xu Jia, Tianyu Zhang, Mengyu Ai, Mang Wang, Yujing Qiao, Lei Zhang, Yanjun Shen, Fan Yang, Miao Zhen, Yijie Zhou, Mingyang Chen, Fei Li, Chenzheng Zhu, Keer Lu, Yaqi Zhao, Hao Liang, Youquan Li, Yanzhao Qin, Linzhuang Sun, Jianhua Xu, Haoze Sun, Mingan Lin, Zenan Zhou, and Weipeng Chen.	1560	
1507	Baichuan-omni-1.5 technical report. <i>Preprint</i> , arXiv:2501.15368.		
1508			
1509			
1510	Yadong Li, Haoze Sun, Mingan Lin, Tianpeng Li, Guosheng Dong, Tao Zhang, Bowen Ding, Wei Song, Zhenglin Cheng, Yuqi Huo, Song Chen, Xu Li, Da Pan, Shusen Zhang, Xin Wu, Zheng Liang, Jun Liu, Tao Zhang, Keer Lu, Yaqi Zhao, Yanjun Shen, Fan Yang, Kaicheng Yu, Tao Lin, Jianhua Xu, Zenan Zhou, and Weipeng Chen.	1561	
1511	Baichuan-omni technical report. <i>Preprint</i> , arXiv:2410.08565.		
1512			
1513			
1514			
1515			
1516	Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen.	2023. Evaluating object hallucination in large vision-language models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 292–305. Association for Computational Linguistics.	1562
1517			1563
1518			1564
1519			1565
1520			1566
1521			1567
1522			1568
1523			1569
1524			1570
1525			1571
1526			1572
1527			1573
1528			1574
1529			1575
1530	Yinghao Aaron Li, Cong Han, and Nima Mesgarani.	2022d. Styletts: A style-based generative model for natural and diverse text-to-speech synthesis. <i>CoRR</i> , abs/2205.15439.	1576
1531			1577
1532			1578
1533			1579
1534			1580
1535			1581
1536			1582
1537			1583
1538	Yizhi Li, Ruibin Yuan, Ge Zhang, Yinghao Ma, Xingran Chen, Hanzhi Yin, Chenghao Xiao, Chenghua Lin, Anton Ragni, Emmanouil Benetos, Norbert Gyenge, Roger B. Dannenberg, Ruibo Liu, Wenhao Chen, Gus Xia, Yemin Shi, Wenhao Huang, Zili Wang, Yike Guo, and Jie Fu.	2024d. MERT: acoustic music understanding model with large-scale self-supervised training. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	1584
1539			1585
1540			1586
1541			1587
1542			1588
1543			1589
1544			1590
1545			1591
1546	Yizhi Li, Ge Zhang, Yinghao Ma, Ruibin Yuan, Kang Zhu, Hangyu Guo, Yiming Liang, Jiaheng Liu, Jian Yang, Siwei Wu, Xingwei Qu, Jinjie Shi, Xinyue Zhang, Zhenzhu Yang, Xiangzhou Wang, Zhaoxiang Zhang, Zachary Liu, Emmanouil Benetos, Wenhao	1592	
1547			1593
1548			1594
1549			1595

1565	Huang, and Chenghua Lin. 2024e. <i>Omnibench: Towards the future of universal omni-language models</i> . <i>CoRR</i> , abs/2409.15272.	Proceedings of Machine Learning Research, pages 21450–21474. PMLR.	1622
1566			1623
1567			
1568	Yunxin Li, Shenyuan Jiang, Baotian Hu, Longyue Wang, Wanqi Zhong, Wenhan Luo, Lin Ma, and Min Zhang. 2024f. <i>Uni-moe: Scaling unified multimodal llms with mixture of experts</i> . <i>CoRR</i> , abs/2405.11273.		
1569			
1570			
1571			
1572	Zhaowei Li, Wei Wang, Yiqing Cai, Qi Xu, Pengyu Wang, Dong Zhang, Hang Song, Botian Jiang, Zhida Huang, and Tao Wang. 2024g. <i>Unifiedilm: Enabling unified representation for multi-modal multi-tasks with large language model</i> . <i>CoRR</i> , abs/2408.02503.		
1573			
1574			
1575			
1576			
1577			
1578	Zhaowei Li, Qi Xu, Dong Zhang, Hang Song, Yiqing Cai, Qi Qi, Ran Zhou, Junting Pan, Zefeng Li, Vu Tu, Zhida Huang, and Tao Wang. 2024h. <i>Groundinggpt: Language enhanced multi-modal grounding model</i> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 6657–6678. Association for Computational Linguistics.		
1579			
1580			
1581			
1582			
1583			
1584			
1585			
1586			
1587	Tian Liang, Jing Huang, Ming Kong, Luyuan Chen, and Qiang Zhu. 2024. <i>Querying as prompt: Parameter-efficient learning for multimodal language model</i> . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 26845–26855. IEEE.		
1588			
1589			
1590			
1591			
1592			
1593	Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. 2024. <i>VILA: on pre-training for visual language models</i> . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 26679–26689. IEEE.		
1594			
1595			
1596			
1597			
1598			
1599	Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. <i>Microsoft COCO: common objects in context</i> . In <i>Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V</i> , volume 8693 of <i>Lecture Notes in Computer Science</i> , pages 740–755. Springer.		
1600			
1601			
1602			
1603			
1604			
1605			
1606			
1607	Samuel Lipping, Parthasarathy Sudarsanam, Konstantinos Drossos, and Tuomas Virtanen. 2022. <i>Clotho-aqa: A crowdsourced dataset for audio question answering</i> . In <i>30th European Signal Processing Conference, EUSIPCO 2022, Belgrade, Serbia, August 29 - Sept. 2, 2022</i> , pages 1140–1144. IEEE.		
1608			
1609			
1610			
1611			
1612			
1613	Fangyu Liu, Guy Emerson, and Nigel Collier. 2023a. <i>Visual spatial reasoning</i> . <i>Trans. Assoc. Comput. Linguistics</i> , 11:635–651.		
1614			
1615			
1616	Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo P. Mandic, Wenwu Wang, and Mark D. Plumbley. 2023b. <i>Audioldm: Text-to-audio generation with latent diffusion models</i> . In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 21450–21474. PMLR.		
1617			
1618			
1619			
1620			
1621			
1622	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. <i>Improved baselines with visual instruction tuning</i> . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 26286–26296. IEEE.		
1623			
1624	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023c. <i>Visual instruction tuning</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .		
1625			
1626			
1627			
1628			
1629			
1630	Shansong Liu, Atin Sakkeer Hussain, Qilong Wu, Chen-shuo Sun, and Ying Shan. 2024b. <i>Mumu-llama: Multi-modal music understanding and generation via large language models</i> . <i>CoRR</i> , abs/2412.06660.		
1631			
1632			
1633			
1634			
1635			
1636	Yaofang Liu, Xiaodong Cun, Xuebo Liu, Xintao Wang, Yong Zhang, Haoxin Chen, Yang Liu, Tieyong Zeng, Raymond Chan, and Ying Shan. 2024c. <i>Evalcrafter: Benchmarking and evaluating large video generation models</i> . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 22139–22149. IEEE.		
1637			
1638			
1639			
1640	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. 2024d. <i>Mmbench: Is your multi-modal model an all-around player?</i> In <i>Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part VI</i> , volume 15064 of <i>Lecture Notes in Computer Science</i> , pages 216–233. Springer.		
1641			
1642			
1643			
1644			
1645			
1646			
1647			
1648	Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. 2022a. <i>Video swin transformer</i> . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 3192–3201. IEEE.		
1649			
1650			
1651			
1652			
1653			
1654			
1655			
1656			
1657	Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. 2022b. <i>A convnet for the 2020s</i> . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 11966–11976. IEEE.		
1658			
1659			
1660			
1661			
1662			
1663	Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. 2024a. <i>Unified-io 2: Scaling autoregressive multimodal models with vision, language, audio, and action</i> . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 26429–26445. IEEE.		
1664			
1665			
1666			
1667			
1668			
1669			
1670			
1671			
1672			
1673			
1674			
1675			
1676			

1677	Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.</i>	Xianzheng Ma, Yash Bhalgat, Brandon Smart, Shuai Chen, Xinghui Li, Jian Ding, Jindong Gu, Dave Zhenyu Chen, Songyou Peng, Jia-Wang Bian, Philip H. S. Torr, Marc Pollefeys, Matthias Nießner, Ian D. Reid, Angel X. Chang, Iro Laina, and Victor Adrian Prisacariu. 2024. When llms step into the 3d world: A survey and meta-analysis of 3d tasks via multi-modal large language models. <i>CoRR</i> , abs/2405.10255.	1733
1678			1734
1679			1735
1680			1736
1681			1737
1682			1738
1683			1739
1684			1740
1685			1741
1686	Pan Lu, Liang Qiu, Jiaqi Chen, Tanglin Xia, Yizhou Zhao, Wei Zhang, Zhou Yu, Xiaodan Liang, and Song-Chun Zhu. 2021. Iconqa: A new benchmark for abstract diagram understanding and visual language reasoning. In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual.</i>	Xiaojian Ma, Silong Yong, Zilong Zheng, Qing Li, Yitao Liang, Song-Chun Zhu, and Siyuan Huang. 2023. SQA3D: situated question answering in 3d scenes. In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.</i> OpenReview.net.	1742
1687			1743
1688			1744
1689			1745
1690			1746
1691			1747
1692			
1693			
1694	Yichen Lu, Jiaqi Song, Xuankai Chang, Hengwei Bian, Soumi Maiti, and Shinji Watanabe. 2024b. Syneslm: A unified approach for audio-visual speech recognition and translation via language model and synthetic data. <i>CoRR</i> , abs/2408.00624.	Muhammad Maaz, Hanoona Abdul Rasheed, Salman Khan, and Fahad Khan. 2024. Video-chatgpt: Towards detailed video understanding via large vision and language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 12585–12602. Association for Computational Linguistics.	1748
1695			1749
1696			1750
1697			1751
1698			1752
1699	Mingshuang Luo, Ruibing Hou, Hong Chang, Zimo Liu, Yaowei Wang, and Shiguang Shan. 2024. M³gpt: An advanced multimodal, multitask framework for motion comprehension and generation. <i>CoRR</i> , abs/2405.16273.	Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. 2023. Egoschema: A diagnostic benchmark for very long-form video language understanding. In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.</i>	1753
1700			1754
1701			1755
1702			1756
1703			
1704	Ruipu Luo, Ziwang Zhao, Min Yang, Junwei Dong, Minghui Qiu, Pengcheng Lu, Tao Wang, and Zhongyu Wei. 2023a. Valley: Video assistant with large language model enhanced ability. <i>CoRR</i> , abs/2306.07207.	Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L. Yuille, and Kevin Murphy. 2016. Generation and comprehension of unambiguous object descriptions. In <i>2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016</i> , pages 11–20. IEEE Computer Society.	1757
1705			1758
1706			1759
1707			1760
1708			1761
1709	Run Luo, Ting-En Lin, Haonan Zhang, Yuchuan Wu, Xiong Liu, Min Yang, Yongbin Li, Longze Chen, Jiaming Li, Lei Zhang, Yangyi Chen, Hamid Alinejad-Rokny, and Fei Huang. 2025. Openomni: Large language models pivot zero-shot omnimodal alignment across language with real-time self-aware emotional speech synthesis. <i>Preprint</i> , arXiv:2501.04561.	Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. OK-VQA: A visual question answering benchmark requiring external knowledge. In <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019</i> , pages 3195–3204. Computer Vision Foundation / IEEE.	1762
1710			1763
1711			
1712			
1713			
1714			
1715			
1716	Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. 2023b. Scalable 3d captioning with pre-trained models. In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.</i>	Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. 2021. Dovqqa: A dataset for VQA on document images. In <i>IEEE Winter Conference on Applications of Computer Vision, WACV 2021, Waikoloa, HI, USA, January 3-8, 2021</i> , pages 2199–2208. IEEE.	1764
1717			1765
1718			1766
1719			1767
1720			1768
1721			1769
1722			1770
1723	Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao, Jingen Zhou, and Tieniu Tan. 2023c. Videofusion: Decomposed diffusion models for high-quality video generation. <i>CoRR</i> , abs/2303.08320.	Xinhao Mei, Chutong Meng, Haohe Liu, Qiuqiang Kong, Tom Ko, Chengqi Zhao, Mark D. Plumbley, Yuexian Zou, and Wenwu Wang. 2024. Wav-caps: A chatgpt-assisted weakly-labelled audio captioning dataset for audio-language multimodal research. <i>IEEE ACM Trans. Audio Speech Lang. Process.</i> , 32:3339–3354.	1771
1724			1772
1725			1773
1726			1774
1727			1775
1728	Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. 2023. Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration. <i>CoRR</i> , abs/2306.09093.		1776
1729			1777
1730			1778
1731			1779
1732			

1790	Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen. 2016. TUT database for acoustic scene classification and sound event detection. In <i>24th European Signal Processing Conference, EUSIPCO 2016, Budapest, Hungary, August 29 - September 2, 2016</i> , pages 1128–1132. IEEE.	1845
1791		1846
1792		1847
1793		1848
1794		1849
1795		1850
1796	OpenBMB MiniCPM-o Team. 2025. Minicpm-o 2.6: A gpt-4o level mllm for vision, speech, and multimodal live streaming on your phone. https://github.com/OpenBMB/MiniCPM-o . Accessed: 2025-02-10.	1851
1797		
1798		
1799		
1800	Anand Mishra, Kartik Alahari, and C. V. Jawahar.	1852
1801	2012. <i>Scene text recognition using higher order language priors</i> . In <i>British Machine Vision Conference, BMVC 2012, Surrey, UK, September 3-7, 2012</i> , pages 1–11. BMVA Press.	1853
1802		1854
1803		1855
1804		1856
1805	Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh,	1857
1806	and Anirban Chakraborty. 2019. <i>OCR-VQA: visual question answering by reading text in images</i> . In <i>2019 International Conference on Document Analysis and Recognition, ICDAR 2019, Sydney, Australia, September 20-25, 2019</i> , pages 947–952. IEEE.	1858
1807		1859
1808		1860
1809		1861
1810		
1811	Seungwhan Moon, Andrea Madotto, Zhaojiang Lin,	1862
1812	Alireza Dirafzoon, Aparajita Saraf, Amy Bear-	1863
1813	man, and Babak Damavandi. 2022. <i>IMU2CLIP: multi-modal contrastive learning for IMU motion</i>	1864
1814	<i>sensors from egocentric videos and text</i> . <i>CoRR</i> ,	1865
1815	abs/2210.14395.	1866
1816		1867
1817	Seungwhan Moon, Andrea Madotto, Zhaojiang Lin,	1868
1818	Tushar Nagarajan, Matt Smith, Shashank Jain, Chun-	1869
1819	Fu Yeh, Prakash Murugesan, Peyman Heidari, Yue	1870
1820	Liu, Kavya Srinet, Babak Damavandi, and Anuj Ku-	1871
1821	marr. 2024. <i>Anymal: An efficient and scalable any-</i>	
1822	<i>modality augmented language model</i> . In <i>Proceed-</i>	
1823	<i>ings of the 2024 Conference on Empirical Methods in Natural Language Processing: EMNLP 2024 - In-</i>	
1824	<i>dustry Track, Miami, Florida, USA, November 12-16,</i>	
1825	<i>2024</i> , pages 1314–1332. Association for Compu-	
1826	tational Linguistics.	
1827		
1828	OpenAI. Hello gpt-4o. https://openai.com/index/hello-gpt-4o/ . Accessed: January 3, 2025.	1872
1829		1873
1830	OpenAI. 2023a. Chatgpt. Technical report, OpenAI.	1874
1831		1875
1832	OpenAI. 2023b. <i>GPT-4 technical report</i> . <i>CoRR</i> ,	1876
1833	abs/2303.08774.	1877
1834		1878
1835	OpenGV. 2024. <i>Interomni</i> . Accessed: 2024-07-27.	
1836		
1837		
1838	Maxime Oquab, Timothée Darcet, Théo Moutakanni,	1884
1839	Huy V. Vo, Marc Szafraniec, Vasil Khalidov,	1885
1840	Pierre Fernandez, Daniel Haziza, Francisco Massa,	1886
1841	Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas,	1887
1842	Wojciech Galuba, Russell Howes, Po-Yao Huang,	1888
1843	Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu	1889
1844	Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jégou,	1890
1845	Julien Mairal, Patrick Labatut, Armand Joulin, and	1891
1846	Piotr Bojanowski. 2024. <i>Dinov2: Learning robust</i>	1892
1847	<i>visual features without supervision</i> . <i>Trans. Mach.</i>	1893
1848	<i>Learn. Res.</i> , 2024.	1894
1849		1895
1850		1896
1851		1897
1852	Vicente Ordonez, Girish Kulkarni, and Tamara L. Berg.	
1853	2011. <i>Im2text: Describing images using 1 million captioned photographs</i> . In <i>Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain</i> , pages 1143–1151.	
1854		
1855		
1856		
1857		
1858		
1859		
1860		
1861		
1862	Aitor Ormazabal, Che Zheng, Cyprien de Masson d'Autume, Dani Yogatama, Deyu Fu, Donovan Ong, Eric Chen, Eugenie Lamprecht, Hai Pham, Isaac Ong, Kaloyan Aleksiev, Lei Li, Matthew Henderson, Max Bain, Mikel Artetxe, Nishant Relan, Piotr Padlewski, Qi Liu, Ren Chen, Samuel Phua, Yazheng Yang, Yi Tay, Yuqi Wang, Zhongkai Zhu, and Zhihui Xie. 2024. <i>Reka core, flash, and edge: A series of powerful multimodal language models</i> . <i>CoRR</i> , abs/2404.12387.	
1863		
1864		
1865		
1866		
1867		
1868		
1869		
1870		
1871		
1872	Artemis Panagopoulou, Le Xue, Ning Yu, Junnan Li,	1862
1873	Dongxu Li, Shafiq Joty, Ran Xu, Silvio Savarese,	1863
1874	Caiming Xiong, and Juan Carlos Niebles. 2024. <i>X-instructclip: A framework for aligning image, 3d, audio, video to llms and its emergent cross-modal reasoning</i> . In <i>Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part XLV</i> , volume 15103 of <i>Lecture Notes in Computer Science</i> , pages 177–197. Springer.	1864
1875		1865
1876		1866
1877		1867
1878		1868
1879	Vassil Panayotov, Guoguo Chen, Daniel Povey, and	1869
1880	Sanjeev Khudanpur. 2015. <i>Librispeech: An ASR corpus based on public domain audio books</i> . In <i>2015 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2015, South Brisbane, Queensland, Australia, April 19-24, 2015</i> , pages 5206–5210. IEEE.	1870
1881		1871
1882		
1883		
1884	Haozhou Pang, Tianwei Ding, Lanshan He, Ming Tao,	1872
1885	Lu Zhang, and Qi Gan. 2024. <i>LLM gesticulator: Leveraging large language models for scalable and controllable co-speech gesture synthesis</i> . <i>CoRR</i> , abs/2410.10851.	1873
1886		1874
1887		1875
1888		1876
1889		1877
1890		1878
1891		
1892	Viorica Patraucean, Lucas Smaira, Ankush Gupta, Adrià Recasens, Larisa Markeeva, Dylan Banarse, Skanda Koppula, Joseph Heyward, Mateusz Malinowski, Yi Yang, Carl Doersch, Tatiana Matejovicova, Yury Sulsky, Antoine Miech, Alexandre Fréchette, Hanna Klimczak, Raphael Koster, Junlin Zhang, Stephanie Winkler, Yusuf Aytar, Simon Osindero, Dima Damen, Andrew Zisserman, and João Carreira. 2023. <i>Perception test: A diagnostic benchmark for multimodal video models</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1879
1893		1880
1894		1881
1895		1882
1896		1883
1897		
1898	Zhiliang Peng, Li Dong, Hangbo Bao, Qixiang Ye, and	1884
1899	Furu Wei. 2022. <i>Beit v2: Masked image modeling with vector-quantized visual tokenizers</i> . <i>CoRR</i> , abs/2208.06366.	1885
1900		1886
1901		1887
1902	Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams,	1888
1903	Luc Van Gool, Markus H. Gross, and Alexander	1889

1904	Sorkine-Hornung. 2016. A benchmark dataset and evaluation methodology for video object segmentation. In <i>2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016</i> , pages 724–732. IEEE Computer Society.	1961
1905		1962
1906		1963
1907		1964
1908		1965
1909	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>J. Mach. Learn. Res.</i> , 21:140:1–140:67.	1966
1910	Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. <i>CoRR</i> , abs/2102.12092.	1967
1911		1968
1912		1969
1913	René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. 2021. Vision transformers for dense prediction. In <i>2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021</i> , pages 12159–12168. IEEE.	1970
1914		1971
1915		1972
1916	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy P. Lillicrap, Jean-Baptiste Alayrac, Radu Soricu, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, Ioannis Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew M. Dai, Katie Millican, Ethan Dyer, Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, James Molloy, Jilin Chen, Michael Isard, Paul Barham, Tom Hennigan, Ross McIlroy, Melvin Johnson, Johan Schalkwyk, Eli Collins, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Clemens Meyer, Gregory Thornton, Zhen Yang, Henryk Michalewski, Zaheer Abbas, Nathan Schucher, Ankesh Anand, Richard Ives, James Keeling, Karel Lenc, Salem Haykal, Siamak Shakeri, Pranav Shyam, Aakanksha Chowdhery, Roman Ring, Stephen Spencer, Eren Sezener, and et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>CoRR</i> , abs/2403.05530.	1973
1917		1974
1918		1975
1919		1976
1920		1977
1921		1978
1922		1979
1923		1980
1924		1981
1925		1982
1926		1983
1927		1984
1928		1985
1929		1986
1930		1987
1931		1988
1932		1989
1933		1990
1934		1991
1935		1992
1936		1993
1937		1994
1938	Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel Synnaeve, and Ronan Collobert. 2020. MLS: A large-scale multilingual dataset for speech research. In <i>21st Annual Conference of the International Speech Communication Association, Interspeech 2020, Virtual Event, Shanghai, China, October 25-29, 2020</i> , pages 2757–2761. ISCA.	1995
1939		1996
1940		1997
1941		1998
1942		1999
1943	Filip Radenovic, Abhimanyu Dubey, Abhishek Kadian, Todor Mihaylov, Simon Vandenhende, Yash Patel, Yi Wen, Vignesh Ramanathan, and Dhruv Mahajan. 2023. Filtering, distillation, and hard negatives for vision-language pre-training. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023</i> , pages 6967–6977. IEEE.	2000
1944		2001
1945	Anhar Risnumawan, Palaiahnakote Shivakumara, Chee Seng Chan, and Chew Lim Tan. 2014. A robust arbitrary text detection system for natural scene images. <i>Expert Syst. Appl.</i> , 41(18):8027–8048.	2002
1946		2003
1947		2004
1948		2005
1949	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 10674–10685. IEEE.	2006
1950		2007
1951		2008
1952		2009
1953		2010
1954	Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2023. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023</i> , pages 22500–22510. IEEE.	2011
1955		2012
1956		2013
1957		2014
1958		2015
1959		2016
1960		2017
1961		2018
1962		2019

2020	RVC-Boss. Gpt-sovits. https://github.com/RVC-Boss/GPT-SoVITS . Accessed: January 3, 2025.	Zhan Shi, Xu Zhou, Xipeng Qiu, and Xiaodan Zhu. 2020. Improving image captioning with better use of captions. <i>CoRR</i> , abs/2006.11807.	2075 2076 2077
2021	Ramon Sanabria, Ozan Caglayan, Shruti Palaskar, Desmond Elliott, Loïc Barrault, Lucia Specia, and Florian Metze. 2018. How2: A large-scale dataset for multimodal language understanding. <i>CoRR</i> , abs/1811.00347.	Fangxun Shu, Lei Zhang, Hao Jiang, and Cihang Xie. 2023. Audio-visual LLM for video understanding. <i>CoRR</i> , abs/2312.06720.	2078 2079 2080
2022	Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. 2022. LAION-5B: an open large-scale dataset for training next generation image-text models. In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	Mustafa Shukor, Corentin Dancette, and Matthieu Cord. 2023. ep-alm: Efficient perceptual augmentation of language models. In <i>IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023</i> , pages 21999–22012. IEEE.	2081 2082 2083 2084 2085 2086
2023	Christoph Schuhmann, Andreas Köpf, Richard Vencu, Theo Coombes, and Romain Beaumont. 2022b(a). Laion-aesthetics.	Gunnar A. Sigurdsson, Gü̈l Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. 2016. Hollywood in homes: Crowd sourcing data collection for activity understanding. In <i>Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I</i> , volume 9905 of <i>Lecture Notes in Computer Science</i> , pages 510–526. Springer.	2087 2088 2089 2090 2091 2092 2093 2094
2024	Christoph Schuhmann, Andreas Köpf, Richard Vencu, Theo Coombes, and Romain Beaumont. 2022b(b). Laion coco: 600m synthetic captions from laion2b-en.	Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. 2012. Indoor segmentation and support inference from RGBD images. In <i>Computer Vision - ECCV 2012 - 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part V</i> , volume 7576 of <i>Lecture Notes in Computer Science</i> , pages 746–760. Springer.	2095 2096 2097 2098 2099 2100 2101
2025	Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuaki. 2021. LAION-400M: open dataset of clip-filtered 400 million image-text pairs. <i>CoRR</i> , abs/2111.02114.	Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards VQA models that can read. In <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019</i> , pages 8317–8326. Computer Vision Foundation / IEEE.	2102 2103 2104 2105 2106 2107 2108
2026	Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022. A-OKVQA: A benchmark for visual question answering using world knowledge. In <i>Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part VIII</i> , volume 13668 of <i>Lecture Notes in Computer Science</i> , pages 146–162. Springer.	Hubert Siuzdak, Florian Grötschla, and Luca A. Lanzendorfer. 2024. SNAC: multi-scale neural audio codec. <i>CoRR</i> , abs/2410.14411.	2109 2110 2111
2027	Share. 2024. Sharegemini: Scaling up video caption data for multimodal large language models.	Shuran Song, Samuel P. Lichtenberg, and Jianxiong Xiao. 2015. SUN RGB-D: A RGB-D scene understanding benchmark suite. In <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015</i> , pages 567–576. IEEE Computer Society.	2112 2113 2114 2115 2116 2117
2028	Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers</i> , pages 2556–2565. Association for Computational Linguistics.	Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. 2012. UCF101: A dataset of 101 human actions classes from videos in the wild. <i>CoRR</i> , abs/1212.0402.	2118 2119 2120 2121
2029	Yiming Shi, Xun Zhu, Ying Hu, Chenyi Guo, Miao Li, and Ji Wu. 2024. Med-2e3: A 2d-enhanced 3d medical multimodal large language model. <i>CoRR</i> , abs/2411.12783.	Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. 2023. Pandagpt: One model to instruction-follow them all. <i>CoRR</i> , abs/2305.16355.	2122 2123 2124
2030		Guangzhi Sun, Wenyi Yu, Changli Tang, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, Yuxuan Wang, and Chao Zhang. 2024a. video-salmonn: Speech-enhanced audio-visual large language models. In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	2125 2126 2127 2128 2129 2130 2131
2031			
2032			
2033			
2034			
2035			
2036			
2037			
2038			
2039			
2040			
2041			
2042			
2043			
2044			
2045			
2046			
2047			
2048			
2049			
2050			
2051			
2052			
2053			
2054			
2055			
2056			
2057			
2058			
2059			
2060			
2061			
2062			
2063			
2064			
2065			
2066			
2067			
2068			
2069			
2070			
2071			
2072			
2073			
2074			

2132	Guangzhi Sun, Wenyi Yu, Changli Tang, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and Chao Zhang. 2023a. <i>Fine-grained audio-visual joint representations for multimodal large language models</i> . <i>CoRR</i> , abs/2310.05863.	2190
2133	Zineng Tang, Ziyi Yang, Mahmoud Khademi, Yang Liu, Chenguang Zhu, and Mohit Bansal. 2024b. <i>Codi-2: In-context, interleaved, and interactive any-to-any generation</i> . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 27415–27424. IEEE.	2191
2134	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2192
2135	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2193
2136	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2194
2137	Keqiang Sun, Junting Pan, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, Jifeng Dai, Yu Qiao, Limin Wang, and Hongsheng Li. 2023b. <i>Journeydb: A benchmark for generative image understanding</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2195
2138	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2196
2139	Yapeng Tian, Dingzeyu Li, and Chenliang Xu. 2020. <i>Unified multisensory perception: Weakly-supervised audio-visual video parsing</i> . In <i>Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part III</i> , volume 12348 of <i>Lecture Notes in Computer Science</i> , pages 436–454. Springer.	2204
2140	Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. 2022. <i>Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training</i> . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	2205
2141	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. <i>Llama: Open and efficient foundation language models</i> . <i>CoRR</i> , abs/2302.13971.	2206
2142	Jack Urbanek, Florian Bordes, Pietro Astolfi, Mary Williamson, Vasu Sharma, and Adriana Romero-Soriano. 2024. <i>A picture is worth more than 77 text tokens: Evaluating clip-style models on dense captions</i> . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 26690–26699. IEEE.	2207
2143	Christophe Veaux, Junichi Yamagishi, Kirsten MacDonald, et al. 2017. <i>Cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit</i> . <i>CSTR</i> , 6:15.	2208
2144	Andreas Veit, Tomas Matera, Lukás Neumann, Jiri Matas, and Serge J. Belongie. 2016. <i>Coco-text: Dataset and benchmark for text detection and recognition in natural images</i> . <i>CoRR</i> , abs/1601.07140.	2209
2145	Jiaqi Wang, Hanqi Jiang, Yiheng Liu, Chong Ma, Xu Zhang, Yi Pan, Mengyuan Liu, Peiran Gu, Sichen Xia, Wenjun Li, Yutong Zhang, Zihao Wu, Zhengliang Liu, Tianyang Zhong, Bao Ge, Tuo Zhang, Ning Qiang, Xintao Hu, Xi Jiang, Xin Zhang, Wei Zhang, Dinggang Shen, Tianming Liu, and Shu	2210
2146	Licai Sun, Zheng Lian, Bin Liu, and Jianhua Tao. 2023c. <i>MAE-DFER: efficient masked autoencoder for self-supervised dynamic facial expression recognition</i> . In <i>Proceedings of the 31st ACM International Conference on Multimedia, MM 2023, Ottawa, ON, Canada, 29 October 2023- 3 November 2023</i> , pages 6110–6121. ACM.	2204
2147	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2205
2148	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2206
2149	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2207
2150	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2208
2151	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2209
2152	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2210
2153	Licai Sun, Zheng Lian, Bin Liu, and Jianhua Tao. 2023d. <i>MAE-DFER: efficient masked autoencoder for self-supervised dynamic facial expression recognition</i> . In <i>Proceedings of the 31st ACM International Conference on Multimedia, MM 2023, Ottawa, ON, Canada, 29 October 2023- 3 November 2023</i> , pages 6110–6121. ACM.	2211
2154	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2212
2155	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2213
2156	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2214
2157	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2215
2158	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2216
2159	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2217
2160	Luoyi Sun, Xuenan Xu, Mengyue Wu, and Weidi Xie. 2024b. <i>Auto-acd: A large-scale dataset for audio-language representation learning</i> . In <i>Proceedings of the 32nd ACM International Conference on Multimedia, MM 2024, Melbourne, VIC, Australia, 28 October 2024 - 1 November 2024</i> , pages 5025–5034. ACM.	2218
2161	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2219
2162	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2220
2163	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2221
2164	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2222
2165	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2223
2166	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2224
2167	Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. 2023e. <i>EVA-CLIP: improved training techniques for CLIP at scale</i> . <i>CoRR</i> , abs/2303.15389.	2225
2168	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2226
2169	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2227
2170	Quan Sun, Qiying Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. 2024c. <i>Emu: Generative pretraining in multimodality</i> . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	2228
2171	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2229
2172	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2230
2173	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2231
2174	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2232
2175	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2233
2176	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2234
2177	Andrew Szot, Bogdan Mazoure, Harsh Agrawal, R. Devon Hjelm, Zsolt Kira, and Alexander Toshev. 2024a. <i>Grounding multimodal large language models in actions</i> . <i>CoRR</i> , abs/2406.07904.	2235
2178	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2236
2179	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2237
2180	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2238
2181	Andrew Szot, Bogdan Mazoure, Omar Attia, Aleksei Timofeev, Harsh Agrawal, Devon Hjelm, Zhe Gan, Zsolt Kira, and Alexander Toshev. 2024b. <i>From multimodal llms to generalist embodied agents: Methods and lessons</i> . <i>Preprint</i> , arXiv:2412.08442.	2239
2182	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2240
2183	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2241
2184	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2242
2185	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2243
2186	Yunlong Tang, Daiki Shimada, Jing Bi, and Chenliang Xu. 2024a. <i>Avicuna: Audio-visual LLM with interleaver and context-boundary alignment for temporal referential dialogue</i> . <i>CoRR</i> , abs/2403.16276.	2244
2187	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2245
2188	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2246
2189	Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. 2023. <i>Any-to-any generation via composable diffusion</i> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2247

2247	Zhang. 2024a. A comprehensive review of multi-modal large language models: Performance and challenges across different tasks. <i>CoRR</i> , abs/2408.01319.	2305
2248		2306
2249		2307
2250	Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. 2023a. <i>Modelscope text-to-video technical report</i> . <i>CoRR</i> , abs/2308.06571.	2308
2251		2309
2252		2310
2253		2311
2254	Kai Wang, Boris Babenko, and Serge J. Belongie. 2011. <i>End-to-end scene text recognition</i> . In <i>IEEE International Conference on Computer Vision, ICCV 2011, Barcelona, Spain, November 6-13, 2011</i> , pages 1457–1464. IEEE Computer Society.	2312
2255		2313
2256		2314
2257		2315
2258		2316
2259	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024b. <i>Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution</i> . <i>CoRR</i> , abs/2409.12191.	2317
2260		2318
2261		2319
2262		2320
2263		2321
2264		2322
2265		
2266		
2267	Wenhai Wang, Jiangwei Xie, Chuanyang Hu, Haoming Zou, Jianan Fan, Wenwen Tong, Yang Wen, Silei Wu, Hanming Deng, Zhiqi Li, Hao Tian, Lewei Lu, Xizhou Zhu, Xiaogang Wang, Yu Qiao, and Jifeng Dai. 2023b. <i>Drivemlm: Aligning multi-modal large language models with behavioral planning states for autonomous driving</i> . <i>CoRR</i> , abs/2312.09245.	2323
2268		2324
2269		2325
2270		2326
2271		2327
2272		2328
2273		
2274	Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. 2019. <i>Vatex: A large-scale, high-quality multilingual dataset for video-and-language research</i> . In <i>2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019</i> , pages 4580–4590. IEEE.	2329
2275		2330
2276		2331
2277		2332
2278		2333
2279		2334
2280		2335
2281	Xinyu Wang, Bohan Zhuang, and Qi Wu. 2024c. <i>Modaverse: Efficiently transforming modalities with llms</i> . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 26596–26606. IEEE.	2336
2282		
2283		
2284		
2285		
2286	Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan Chen, Yaohui Wang, Ping Luo, Ziwei Liu, Yali Wang, Limin Wang, and Yu Qiao. 2024d. <i>Internvid: A large-scale video-text dataset for multimodal understanding and generation</i> . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	2337
2287		2338
2288		2339
2289		2340
2290		2341
2291		
2292		
2293		
2294	Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun Zheng, Zun Wang, Yansong Shi, Tianxiang Jiang, Songze Li, Jian Xu, Hongjie Zhang, Yifei Huang, Yu Qiao, Yali Wang, and Limin Wang. 2024e. <i>Internvideo2: Scaling foundation models for multimodal video understanding</i> . In <i>Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXXXV</i> , volume 15143 of <i>Lecture Notes in Computer Science</i> , pages 396–416. Springer.	2342
2295		2343
2296		2344
2297		2345
2298		2346
2299		
2300		
2301		
2302		
2303		
2304		
2305	Zekun Wang, King Zhu, Chunpu Xu, Wangchunshu Zhou, Jiaheng Liu, Yibo Zhang, Jiashuo Wang, Ning Shi, Siyu Li, Yizhi Li, Haoran Que, Zhaoxiang Zhang, Yuanxing Zhang, Ge Zhang, Ke Xu, Jie Fu, and Wenhao Huang. 2024f. <i>Mio: A foundation model on multimodal tokens</i> . <i>Preprint</i> , arXiv:2409.17692.	2347
2306		2348
2307		2349
2308		2350
2309		2351
2310		2352
2311		2353
2312	Julong Wei, Shanshuai Yuan, Pengfei Li, Qingda Hu, Zhongxue Gan, and Wenchao Ding. 2024. <i>Occlama: An occupancy-language-action generative world model for autonomous driving</i> . <i>CoRR</i> , abs/2409.03272.	2313
2313		2314
2314		2315
2315		2316
2316		
2317	Bo Wu, Shoubin Yu, Zhenfang Chen, Josh Tenenbaum, and Chuang Gan. 2021. <i>STAR: A benchmark for situated reasoning in real-world videos</i> . In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .	2317
2318		2318
2319		2319
2320		2320
2321		2321
2322		2322
2323	Jiahong Wu, He Zheng, Bo Zhao, Yixin Li, Baoming Yan, Rui Liang, Wenjia Wang, Shipei Zhou, Guosen Lin, Yanwei Fu, Yizhou Wang, and Yonggang Wang. 2017. <i>AI challenger : A large-scale dataset for going deeper in image understanding</i> . <i>CoRR</i> , abs/1711.06475.	2323
2324		2324
2325		2325
2326		2326
2327		2327
2328		2328
2329	Jialian Wu, Jianfeng Wang, Zhengyuan Yang, Zhe Gan, Zicheng Liu, Junsong Yuan, and Lijuan Wang. 2024a. <i>Grit: A generative region-to-text transformer for object understanding</i> . In <i>Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXXX, volume 15138 of Lecture Notes in Computer Science</i> , pages 207–224. Springer.	2329
2330		2330
2331		2331
2332		2332
2333		2333
2334		2334
2335		2335
2336		2336
2337	Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and Philip S. Yu. 2023. <i>Multimodal large language models: A survey</i> . In <i>IEEE International Conference on Big Data, BigData 2023, Sorrento, Italy, December 15-18, 2023</i> , pages 2247–2256. IEEE.	2337
2338		2338
2339		2339
2340		2340
2341		2341
2342	Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2024b. <i>Next-gpt: Any-to-any multimodal LLM</i> . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	2342
2343		2343
2344		2344
2345		2345
2346		2346
2347	Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 2015. <i>3d shapenets: A deep representation for volumetric shapes</i> . In <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015</i> , pages 1912–1920. IEEE Computer Society.	2347
2348		2348
2349		2349
2350		2350
2351		2351
2352		2352
2353		2353
2354	Hanguang Xiao, Feizhong Zhou, Xingyue Liu, Tianqi Liu, Zhipeng Li, Xin Liu, and Xiaoxuan Huang. 2024. <i>A comprehensive survey of large language models and multimodal large language models in medicine</i> . <i>CoRR</i> , abs/2405.08603.	2354
2355		2355
2356		2356
2357		2357
2358		2358
2359	Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. <i>Next-qa: Next phase of question-answering to explaining temporal actions</i> . In <i>IEEE</i>	2359
2360		2360
2361		2361

Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021, pages 9777–9786. Computer Vision Foundation / IEEE.

Zhifei Xie and Changqiao Wu. 2024. Mini-omni2: Towards open-source gpt-4o with vision, speech and duplex capabilities. *CoRR*, abs/2410.11190.

Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueteng Zhuang. 2017. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 2017 ACM on Multimedia Conference, MM 2017, Mountain View, CA, USA, October 23-27, 2017*, pages 1645–1653. ACM.

Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofighi, Fisher Yu, Dacheng Tao, and Andreas Geiger. 2023. Unifying flow, stereo and depth estimation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 45(11):13941–13958.

Jingwei Xu, Chenyu Wang, Zibo Zhao, Wen Liu, Yi Ma, and Shenghua Gao. 2024a. CAD-MLLM: unifying multimodality-conditioned CAD generation with MLLM. *CoRR*, abs/2411.04954.

Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. MSR-VTT: A large video description dataset for bridging video and language. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 5288–5296. IEEE Computer Society.

Runsen Xu, Xiaolong Wang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. 2024b. Pointllm: Empowering large language models to understand point clouds. In *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part XXV*, volume 15083 of *Lecture Notes in Computer Science*, pages 131–147. Springer.

Jinlong Xue, Yayue Deng, Yingming Gao, and Ya Li. 2024a. Auffusion: Leveraging the power of diffusion and large language models for text-to-audio generation. *IEEE ACM Trans. Audio Speech Lang. Process.*, 32:4700–4712.

Le Xue, Ning Yu, Shu Zhang, Artemis Panagopoulou, Junnan Li, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, and Silvio Savarese. 2024b. ULIP-2: towards scalable multimodal pre-training for 3d understanding. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 27081–27091. IEEE.

Honghui Yang, Tong He, Jiaheng Liu, Hua Chen, Boxi Wu, Binbin Lin, Xiaofei He, and Wanli Ouyang. 2023a. GD-MAE: generative decoder for MAE pre-training on lidar point clouds. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 9403–9414. IEEE.

Zhen Yang, Yingxue Zhang, Fandong Meng, and Jie Zhou. 2023b. TEAL: tokenize and embed ALL for multi-modal large language models. *CoRR*, abs/2311.04589.

Hanrong Ye, De-An Huang, Yao Lu, Zhiding Yu, Wei Ping, Andrew Tao, Jan Kautz, Song Han, Dan Xu, Pavlo Molchanov, and Hongxu Yin. 2024a. X-VILA: cross-modality alignment for large language model. *CoRR*, abs/2405.19335.

Qilang Ye, Zitong Yu, Rui Shao, Xinyu Xie, Philip Torr, and Xiaochun Cao. 2024b. CAT: enhancing multi-modal large language model to answer questions in dynamic audio-visual scenarios. In *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part X*, volume 15068 of *Lecture Notes in Computer Science*, pages 146–164. Springer.

Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2023a. A survey on multimodal large language models. *CoRR*, abs/2306.13549.

Zhenfei Yin, Jiong Wang, Jianjian Cao, Zhelun Shi, Dingning Liu, Mukai Li, Xiaoshui Huang, Zhiyong Wang, Lu Sheng, Lei Bai, Jing Shao, and Wanli Ouyang. 2023b. LAMM: language-assisted multi-modal instruction-tuning dataset, framework, and benchmark. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. 2024. Yi: Open foundation models by 01.ai. *CoRR*, abs/2403.04652.

Jiazu Yu, Haomiao Xiong, Lu Zhang, Haiwen Diao, Yunzhi Zhuge, Lanqing Hong, Dong Wang, Huchuan Lu, You He, and Long Chen. 2024a. Llms can evolve continually on modality for x-modal reasoning. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.

Licheng Yu, Patrick Poirson, Shan Yang, Alexander C. Berg, and Tamara L. Berg. 2016. Modeling context in referring expressions. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II*, volume 9906 of *Lecture Notes in Computer Science*, pages 69–85. Springer.

Lijun Yu, José Lezama, Nitesh Bharadwaj Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong Cheng, Agrim Gupta, Xiuye Gu, Alexander G. Hauptmann, Boqing Gong, Ming-Hsuan Yang, Irfan Essa,

2476	David A. Ross, and Lu Jiang. 2024b. Language model beats diffusion - tokenizer is key to visual generation. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	2534
2477		2535
2478		2536
2479		2537
2480		2538
2481	Shoubin Yu, Jaehong Yoon, and Mohit Bansal. 2024c. Crema: Generalizable and efficient video-language reasoning via multimodal modular fusion. <i>Preprint</i> , arXiv:2402.05889.	2539
2482		2540
2483		2541
2484		2542
2485	Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2024d. Mm-vet: Evaluating large multimodal models for integrated capabilities. In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	2543
2486		2544
2487		2545
2488		2546
2489		2547
2490		2548
2491		2549
2492	Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yuet-ing Zhuang, and Dacheng Tao. 2019. Activitynet-qa: A dataset for understanding complex web videos via question answering. In <i>The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019</i> , pages 9127–9134. AAAI Press.	2550
2493		2551
2494		2552
2495		2553
2496		2554
2497		2555
2498		2556
2499		2557
2500		2558
2501		2559
2502		
2503	Xiang Yue, Yuansheng Ni, Tianyu Zheng, Kai Zhang, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. 2024. MMMU: A massive multi-discipline multimodal understanding and reasoning benchmark for expert AGI. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 9556–9567. IEEE.	2560
2504		2561
2505		2562
2506		2563
2507		2564
2508		2565
2509		2566
2510		2567
2511		
2512		
2513		
2514	Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. 2022. Sound-stream: An end-to-end neural audio codec. <i>IEEE ACM Trans. Audio Speech Lang. Process.</i> , 30:495–507.	2568
2515		2569
2516		2570
2517		2571
2518		2572
2519	Rowan Zellers, Jiasen Lu, Ximing Lu, Youngjae Yu, Yanpeng Zhao, Mohammadreza Salehi, Aditya Kusupati, Jack Hessel, Ali Farhadi, and Yejin Choi. 2022. MERLOT RESERVE: neural script knowledge through vision and language and sound. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 16354–16366. IEEE.	2573
2520		2574
2521		2575
2522		
2523		
2524		
2525		
2526		
2527	Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang,	2576
2528		2577
2529		2578
2530		2579
2531		2580
2532		2581
2533		2582
510	Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. 2024. Chatglm: A family of large language models from GLM-130B to GLM-4 all tools. <i>CoRR</i> , abs/2406.12793.	2583
511		2584
512		2585
513		2586
514		2587
515		2588
516		2589
517		2590
518		
519	Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. In <i>IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023</i> , pages 11941–11952. IEEE.	2544
520		2545
521		2546
522		2547
523		2548
524		2549
525		
526	Jun Zhan, Junqi Dai, Jiasheng Ye, Yunhua Zhou, Dong Zhang, Zhigeng Liu, Xin Zhang, Ruibin Yuan, Ge Zhang, Linyang Li, Hang Yan, Jie Fu, Tao Gui, Tianxiang Sun, Yu-Gang Jiang, and Xipeng Qiu. 2024. Anygpt: Unified multimodal LLM with discrete sequence modeling. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 9637–9662. Association for Computational Linguistics.	2550
527		2551
528		2552
529		2553
530		2554
531		2555
532		2556
533		2557
534		2558
535		2559
536		
537	Binbin Zhang, Hang Lv, Pengcheng Guo, Qijie Shao, Chao Yang, Lei Xie, Xin Xu, Hui Bu, Xiaoyu Chen, Chenchen Zeng, Di Wu, and Zhendong Peng. 2022a. WENETSPEECH: A 10000+ hours multi-domain mandarin corpus for speech recognition. In <i>IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022</i> , pages 6182–6186. IEEE.	2560
538		2561
539		2562
540		2563
541		2564
542		2565
543		2566
544		2567
545		
546	Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. 2023a. Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 15757–15773. Association for Computational Linguistics.	2568
547		2569
548		2570
549		2571
550		2572
551		2573
552		2574
553		2575
554		
555	Duzhen Zhang, Yahan Yu, Jiahua Dong, Chenxing Li, Dan Su, Chenhui Chu, and Dong Yu. 2024a. Mmlms: Recent advances in multimodal large language models. In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 12401–12430. Association for Computational Linguistics.	2576
556		2577
557		2578
558		2579
559		2580
560		2581
561		2582
562		
563	Hang Zhang, Xin Li, and Lidong Bing. 2023b. Video-llama: An instruction-tuned audio-visual language model for video understanding. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023 - System Demonstrations, Singapore, December 6-10, 2023</i> , pages 543–553. Association for Computational Linguistics.	2583
564		2584
565		2585
566		2586
567		2587
568		2588
569		2589
570		2590

- 2591 Meishan Zhang, Hao Fei, Bin Wang, Shengqiong Wu,
 2592 Yixin Cao, Fei Li, and Min Zhang. 2024b. Recogniz-
 2593 ing everything from all modalities at once: Grounded
 2594 multimodal universal information extraction. In *Find-
 2595 ings of the Association for Computational Linguistics,
 2596 ACL 2024, Bangkok, Thailand and virtual meeting,
 2597 August 11-16, 2024*, pages 14498–14511. Associa-
 2598 tion for Computational Linguistics.
- 2599 Susan Zhang, Stephen Roller, Naman Goyal, Mikel
 2600 Artetxe, Moya Chen, Shuhui Chen, Christopher
 2601 Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin,
 2602 Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shus-
 2603 ter, Daniel Simig, Punit Singh Koura, Anjali Srid-
 2604 har, Tianlu Wang, and Luke Zettlemoyer. 2022b.
 2605 **OPT: open pre-trained transformer language mod-
 2606 els.** *CoRR*, abs/2205.01068.
- 2607 Xin Zhang, Dong Zhang, Shimin Li, Yaqian Zhou, and
 2608 Xipeng Qiu. 2023c. **Speecktokenizer: Unified speech
 2609 tokenizer for speech large language models.** *CoRR*,
 2610 abs/2308.16692.
- 2611 Yiyuan Zhang, Kaixiong Gong, Kaipeng Zhang, Hong-
 2612 sheng Li, Yu Qiao, Wanli Ouyang, and Xiangyu Yue.
 2613 2023d. **Meta-transformer: A unified framework for
 2614 multimodal learning.** *CoRR*, abs/2307.10802.
- 2615 Yang Zhao, Zhijie Lin, Daquan Zhou, Zilong Huang,
 2616 Jiashi Feng, and Bingyi Kang. 2023a. **Bubogpt: En-
 2617 abling visual grounding in multi-modal llms.** *CoRR*,
 2618 abs/2307.08581.
- 2619 Zijia Zhao, Longteng Guo, Tongtian Yue, Sihan Chen,
 2620 Shuai Shao, Xinxin Zhu, Zehuan Yuan, and Jing
 2621 Liu. 2023b. **Chatbridge: Bridging modalities with
 2622 large language model as a language catalyst.** *CoRR*,
 2623 abs/2305.16103.
- 2624 Zhisheng Zhong, Chengyao Wang, Yuqi Liu, Senqiao
 2625 Yang, Longxiang Tang, Yuechen Zhang, Jingyao Li,
 2626 Tianyuan Qu, Yanwei Li, Yukang Chen, Shaozuo Yu,
 2627 Sitong Wu, Eric Lo, Shu Liu, and Jiaya Jia. 2024.
 2628 **Lyra: An efficient and speech-centric framework for
 2629 omni-cognition.** *Preprint*, arXiv:2412.09501.
- 2630 Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiaxi Cui,
 2631 Hongfa Wang, Yatian Pang, Wenhao Jiang, Junwu
 2632 Zhang, Zongwei Li, Caiwan Zhang, Zhipeng Li,
 2633 Wei Liu, and Li Yuan. 2024a. **Languagebind: Ex-
 2634 tends video-language pretraining to n-modality by
 2635 language-based semantic alignment.** In *The Twelfth
 2636 International Conference on Learning Representa-
 2637 tions, ICLR 2024, Vienna, Austria, May 7-11, 2024.*
 2638 OpenReview.net.
- 2639 Hongyan Zhu, Shuai Qin, Min Su, Chengzhi Lin, Anjie
 2640 Li, and Junfeng Gao. 2024b. **Harnessing large vision
 2641 and language models in agriculture: A review.** *CoRR*,
 2642 abs/2407.19679.
- 2643 Wanrong Zhu, Jack Hessel, Anas Awadalla,
 2644 Samir Yitzhak Gadre, Jesse Dodge, Alex Fang,
 2645 Youngjae Yu, Ludwig Schmidt, William Yang Wang,
 2646 and Yejin Choi. 2023a. **Multimodal C4: an open,
 2647 billion-scale corpus of images interleaved with**
- text. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.*
- 2648
 2649
 2650
 2651
 2652
 2653
 2654
 2655
 2656
 2657
 2658
 2659
 2660

2661 A Related Survey

2662 With the advent of MLLMs, there are several surveys detailing the current progress of MLLMs.
2663 Yin et al. (2023a); Wu et al. (2023); Caffagni
2664 et al. (2024) focus on the early Vision-MLLMs,
2665 while Çoban et al. (2024) and Ma et al. (2024)
2666 respectively summarize the Audio-MLLMs and
2667 3D-MLLMs. Zhang et al. (2024a); Wang et al.
2668 (2024a) conduct an investigation into various
2669 Specific-MLLMs of different modalities. He et al.
2670 (2024); Chen et al. (2024b) discuss the expansion
2671 of MLLM’s generative capabilities. Some
2672 works discuss MLLMs in specific domains, such
2673 as medicine (Xiao et al., 2024), agriculture (Zhu
2674 et al., 2024b), and autonomous driving (Cui et al.,
2675 2024). Some works highlight some specific tasks
2676 such as safety (Fan et al., 2024), hallucination (Bai
2677 et al., 2024b), and acceleration (Zhu et al., 2024b).
2678 And Li and Lu (2024); Huang and Zhang (2024)
2679 focus on the evaluation of MLLM performance.

2680 Distinct from the above-mentioned surveys,
2681 this paper focuses on MLLMs that align multiple
2682 non-linguistic modalities² with LLMs (Omni-
2683 MLLMs), enabling cross-modal understanding or
2684 cross-modal generation. As the first systematic
2685 survey on Omni-MLLMs, we hope our work will
2686 serve as an overview of this emerging direction,
2687 fostering future research in the field.

2689 B Details about Omni-MLLMs architectures

2690 Table 1 presents the details of the structure of main-
2691 stream Omni-MLLMs. We will list some of the
2692 pre-trained models used.

2693 B.1 Modality Encoder

2694 **Visual Specific-Encoder** Vit (Dosovitskiy et al.,
2695 2021), SigCLIP Vit (Zhai et al., 2023), CLIP
2696 Vit (Radford et al., 2021), EVA CLIP Vit (Sun et al.,
2697 2023e), InternVit (Chen et al., 2023g), DINOV2
2698 Vit (Oquab et al., 2024), DFNCLIP Vit (Fang
2699 et al., 2024a), and OpenCLIP ConvNext (Liu
2700 et al., 2022b) encode images to obtain continu-
2701 ous features. TimeSformer (Bertasius et al., 2021),
2702 VideoMAE (Tong et al., 2022), MAE-DFFer (Sun
2703 et al., 2023c), Omni-VL (Sun et al., 2023d), Video-
2704 Swin (Liu et al., 2022a), and Vivit (Arnab et al.,
2705 2021) encode videos to obtain continuous features.

²Since MLLMs capable of comprehending both video and imagery generally process video as multiple frames and employ a single vision encoder, we categorize them as Specific-MLLMs, i.e. Vision-MLLMs.

2706 **Audio Specific-Encoder** AST (Gong et al.,
2707 2021), Beats (Chen et al., 2023d), Whisper (Rad-
2708 ford et al., 2023), HuBERT (Hsu et al., 2021),
2709 CLAP (Elizalde et al., 2023), Conformer (Gu-
2710 lati et al., 2020), MERT (Li et al., 2024d), and
2711 PANN (Kong et al., 2020) encode the audio modal-
2712 ity to obtain continuous features.

2713 **3D Specific-Encoder** ULIP2 (Xue et al., 2024b),
2714 GD-MAE (Yang et al., 2023a), PointEncoder (Xu
2715 et al., 2024b), FrozenCLIP (Huang et al., 2022),
2716 and M3D-CLIP (Bai et al., 2024a) encode the 3D
2717 modality to obtain continuous features.

2718 **Pre-align Uni-Encoder** LanguageBind (Zhu
2719 et al., 2024a), ImageBind (Girdhar et al., 2023),
2720 Meta-Transformers (Zhang et al., 2023d), TVL (Fu
2721 et al., 2024c), and SSVTP (Kerr et al., 2023)
2722 encode multiple non-linguistic modalities into a uni-
2723 fied feature space and obtain continuous features.
2724 TVL, LanguageBind, ImageBind, and SSVTP
2725 construct modality-specific encoders for different
2726 modalities and achieve multi-modalities alignment
2727 through indirect alignment. Meta-Transformers de-
2728 sign distinct modality-specific patch embeddings
2729 and use a shared encoder to encode multiple modal-
2730 ities.

2731 **Other Specific-Encoder** IMU2CLIP (Moon
2732 et al., 2022) encodes the IMU modality to obtain
2733 continuous features. Individual modality-specific
2734 encoders from LanguageBind or ImageBind are
2735 often used independently as specific encoders.

2736 B.2 Modality Tokenizer

2737 **Visual Tokenizer** VQ-GAN (Esser et al., 2021),
2738 DALL-E (Ramesh et al., 2021), BEiT-V2 (Peng
2739 et al., 2022), MAGVIT-v2 (Yu et al., 2024b), and
2740 SEED (Ge et al., 2023) encode the visual modality
2741 into discrete visual tokens, which can be decoded
2742 back into the original image using the de-tokenizer.

2743 **Audio Tokenizer** Jukebox (Dhariwal et al.,
2744 2020), SoundStream (Zeghidour et al., 2022),
2745 SpeechTokenizer (Zhang et al., 2023c), En-
2746 codec (Défossez et al., 2023), and S2U (Chen et al.,
2747 2024c) encode the audio modality into discrete au-
2748 dio tokens, which can be decoded back into the
2749 audio using the corresponding de-tokenizer.

2750 **Other Tokenizer** Scene Tokenizer (Wei et al.,
2751 2024) encodes the 3D modality into discrete 3D
2752 tokens. LEO (Huang et al., 2024), Ground-
2753 Action (Szot et al., 2024a), OccLLaMA (Wei et al.,

2754, and GMA (Szot et al., 2024b) perform discrete encoding of the action modality to obtain corresponding action tokens, which can be decoded back into the original action using the corresponding de-tokenizer. M3GPT (Luo et al., 2024), Gesticator (Pang et al., 2024), and SOLAMI (Jiang et al., 2024b) perform discrete encoding of the motion modality to obtain corresponding motion tokens, which can be decoded back into the original motion using the corresponding de-tokenizer.

B.3 Modality Generation Model

For image generation, Stable Diffusion (Rombach et al., 2022) and Instruct-Pix2Pix (Brooks et al., 2023) are used. Video generation models include Zeroscope (Cerspense, 2023), VideoFusion (Luo et al., 2023c), VideoCrafter (Chen et al., 2023b), and ModelScope (Wang et al., 2023a). For audio generation, models such as AudioLDM (Liu et al., 2023b), SNAc (Siuzdak et al., 2024), LLaMA-Omni’s audio decoder (Fang et al., 2024b), MusicGen (Copet et al., 2023), and TiCodec (Ren et al., 2024) are utilized. Meanwhile, StyleTTS (Li et al., 2022d) and GPT-SoVITS (RVC-Boss) are employed for speech generation.

B.4 LLM Backbone

Commonly used LLMs include the T5 series (Rafel et al., 2020), LLaMA series (Touvron et al., 2023), Qwen series (Bai et al., 2023), Internlm series (Cai et al., 2024), Chatglm series (Zeng et al., 2024), OPT series (Zhang et al., 2022b), Mixtral series (Jiang et al., 2024a), Mistral series (Jiang et al., 2023), Phi series (Gunasekar et al., 2023), and Yi series (Young et al., 2024).

C Details of Training and evaluation

C.1 Details of Training Data

The statistical results of some commonly used alignment datasets and the instruction data of mainstream Omni-MLLMs are shown in Table 2 and Table 3. There is still a lack of alignment data for data-scarcity modalities and cross-modal instruction data.

C.2 Details of Benchmark

The statistical data of some commonly used benchmarks are shown in Table 4. Existing benchmarks still require improvements in terms of the number of modalities and the forms of modality interaction.

C.3 Other Train Recipes

In addition to the general training paradigms mentioned in Section 3, some other useful training recipes are also used. (1) **Prior knowledge from Specific-MLLMs**: Since Specific-MLLMs have already achieved effective alignment in single-modal scenarios, some Omni-MLLMs directly leverage their well-trained projectors to reduce the training overhead during the alignment phase. For example, InstructBLIP (Panagopoulou et al., 2024) and X-LLM (Chen et al., 2023a) use the Q-former trained by BLIP2 to align the visual modality, while NaviveMC and DAMC (Chen et al., 2024a) further leverage projectors from multiple models to handle alignment for visual, audio, and 3D modalities separately; (2) **Additional human preference training**: Szot et al. (2024b) and Ye et al. (2024b) adopt HF training methods like PPO and ADPO to better align with human preferences; (3) **Modalities Blending**: During progressive alignment pre-training or multi-step instruction fine-tuning, some works (Han et al., 2024a; Li et al., 2024c; Chen et al., 2024c) mix previously trained modality data with the current new modality data for training to prevent catastrophic forgetting.

C.4 Performance of Omni-MLLMs

We statistic the performance of various mainstream Omni-MLLMs in uni-modal understanding, cross-modal understanding, and cross-modal, as shown in Table 5. We also show the performance of several Specific-MLLMs (Lin et al., 2024; Li et al., 2024a; Chu et al., 2023; Xu et al., 2024b; Sun et al., 2024c; Jin et al., 2024) on selected tasks for comparison. The results are mainly from corresponding papers (some results are used as baselines in other papers). It is worth noting that due to differences in the size and performance of the pre-trained models, Omni-MLLMs with the same backbone LLM may still not be fairly comparable. Therefore, this table only provides a rough trend of performance.

It can be seen from the table that most Omni-MLLMs still exhibit a significant performance gap in uni-modal understanding tasks compared to Specific-MLLMs. Meanwhile, in uni-modal generation tasks, models like AnyGPT and CoDi-2 have achieved performance close to or even surpassing Specific-MLLMs. Additionally, Omni-MLLMs are capable of performing cross-modal tasks that Specific-MLLMs cannot handle.

Model	Capabilities	Modalities	Multi-Modalities Encoding		Method	Multi-Modalities Alignment		Vocabulary	Multi-Modalities Interaction	Modalities	Multi-Modalities Generation	
			Method	Encoding Model		Projector	VLLM				Method	Generation model
e-ALM	Cross-modal Understanding	Visual/Audio	Continuous Encoding	Vit/TimeFormer+AST	multi-branch	linear	-	injection	OPT	-	-	-
VALOR	Cross-modal Understanding	Visual/Audio	Continuous Encoding	Vit/VideoSwin+AST	multi-branch	MLP	-	injection	Bert	-	-	-
X-LLM	Cross-modal Understanding	Visual/Audio	Continuous Encoding	Vit/Conformer	multi-branch	Q-former+Linear	-	concatenate	ChatGLM	-	-	-
CharBridge	Cross-modal Understanding	Visual/Audio	Continuous Encoding	EVAL CLIP Vit/Beats	multi-branch	Precovver	-	concatenate	Vicuna	-	-	-
PandaGPT	Cross-modal Understanding	Visual/Audio/3D/IMU/thermal	Continuous Encoding	CLIP Vision+Audio	multi-branch	linear	-	concatenate	Vicuna	-	-	-
VideoLlama	Cross-modal Understanding	Visual/Audio	Continuous Encoding	EVA CLIP Vit/	multi-branch	Q-former+Linear	-	concatenate	Vicuna	-	-	-
LAMM	Cross-modal Understanding	Visual/Audio	Continuous Encoding	ImageBind	multi-branch	CLIP Vit/	-	concatenate	Vicuna	-	-	-
Macaw-LLM	Cross-modal Understanding	Visual/Audio	Continuous Encoding	FrozenCLIP	multi-branch	MLP	-	concatenate	Vicuna	-	-	-
BuboGPT	Cross-modal Understanding	Visual/Audio	Continuous Encoding	Vit/Whisper	multi-branch	Cross-Attention	-	concatenate	LLaMA	-	-	-
Teal	Cross-modal Understanding & Generation	Visual/Audio	Discrete Encoding	CLIP Vision+Audio	multi-branch	Q-former+Linear	-	concatenate	LLaMA	-	-	-
ImageBind-LLM	Cross-modal Understanding & Generation	Visual/Audio/3D	Continuous Encoding	Whisper+K-means	embedding	-	Added Vocabulary	concatenate	LLaMA	Image	Modality-Token-based	VQGAN denoizer
NextGPT	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	ImageBind	uni-branch	MLP	-	injection	LLaMA	-	-	-
Any-MAL	Cross-modal Understanding	Visual/Audio/TMU	Continuous Encoding	CLIP Vit/CLAP/IMU2CLIP	multi-branch	Precovver	-	concatenate	LLaMA-2	-	-	-
FAVOR	Cross-modal Understanding	Visual/Audio	Continuous Encoding	EVA CLIP Vit/Beats	multi-branch	Q-former+Linear	-	concatenate	Vicuna	-	-	-
Octavia	Cross-modal Understanding	Visual/3D	Continuous Encoding	CLIP Vit/Object-A-Scene	multi-branch	MLP	-	concatenate	Vicuna	-	-	-
LEO	Cross-modal Understanding & Generation	Visual/3D/Action	Hybrid Encoding	OpenCLIP Convnet/PointNet++/LEO's Action tokenizer	multi-branch	Spatial Transformer	Overwrite Vocabulary	concatenate	Vicuna	Action	Modality-Token-based	LEO's Action denoizer
CoDi-2	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	ImageBind	multi-branch	Linear	-	concatenate	LLaMA-2	Image/Video/Audio	Representation-based	StableDiffusion2.1/Zeroscope/AudioLM
X-InstructBLIP	Cross-modal Understanding	Visual/Audio/3D	Continuous Encoding	EVA CLIP Vit/Beats/ULP2	multi-branch	Q-former+Linear	-	concatenate	Vicuna	-	-	-
One-LLM	Cross-modal Understanding	Visual/Audio/3D/IMU/URFU	Continuous Encoding	Meta-transformer	uni-branch	UPM(self-attention)	-	concatenate	LLaMA-2	-	-	-
AV-LLM	Cross-modal Understanding	Visual/Audio	Continuous Encoding	CLIP Vit/CLAP	multi-branch	Linear	-	concatenate	Vicuna	-	-	-
DriveMLM	Cross-modal Understanding	Visual/3D	Continuous Encoding	EVA CLIP Vit/GD-MAE	multi-branch	Q-former+Linear	-	concatenate	LLaMA	-	-	-
Omni-3D	Cross-modal Understanding	Visual/3D	Continuous Encoding	CLIP Vit/PointNet++	multi-branch	MLP	-	concatenate	LLaMA-2	-	-	-
ModVerse	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	ImageBind	multi-branch	Linear	-	concatenate	Vicuna	Image/Video/Audio	Text-based	StableDiffusion /AudioLM/VideoFusion
MuhiPLY	Cross-modal Understanding & Thermal/Touch	Visual/Audio/3D	Continuous Encoding	CLIP Vit/CLAP /ConceptGrasp	multi-branch	MLP/Linear	-	concatenate	Vicuna	-	-	-
CREMA	Cross-modal Understanding & Thermal/Touch	Visual/Audio/3D	Continuous Encoding	EVA CLIP Vit/Lineart/ /Thermal/Touch	uni-branch	Q-former+Linear	-	concatenate	Flan-T5	-	-	-
GroundingGPT	Cross-modal Understanding & Thermal/Touch	Visual/Audio	Continuous Encoding	BeamerCLIP /CLIP Vit/Beats/PointNet++/CLIP Vit/Beats/PointBert/PointLM	multi-branch	Q-former+Linear/MLP	-	concatenate	Vicuna	-	-	-
DAMC	Cross-modal Understanding	Visual/Audio	Continuous Encoding	CLIP Vit/Beats/PointBert/PointLM	multi-branch	Q-former+Linear/MLP	-	concatenate	Vicuna	-	-	-
AnyGPT	Cross-modal Understanding & Generation	Visual/Audio	Discrete Encoding	SEED tokenizer	embedding	-	Extend Vocabulary	concatenate	LLaMA-2	Image/Speech/Music	-	SEED de-tokenizer/Speech de-tokenizer/Encoder-de-tokenizer
TVL-LLaMA	Cross-modal Understanding	Visual/Touch	Continuous Encoding	EVA CLIP Vit/Encoder	uni-branch	MLP	-	injection	LLaMA	-	-	-
SSVT-LLaMA	Cross-modal Understanding	Visual/Touch	Continuous Encoding	SSVT Encoders	uni-branch	MLP	-	injection	LLaMA	-	-	-
CAT	Cross-modal Understanding	Visual/Audio	Continuous Encoding	ImageBind	multi-branch	Linear	-	concatenate	LLaMA-2	-	-	-
Vicuna	Cross-modal Understanding	Visual/Audio	Continuous Encoding	CLIP Vit/CLAP	multi-branch	MLP	-	concatenate	Vicuna	-	-	-
WorldGPT	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	LanguageBind	uni-branch	Linear	-	concatenate	Vicuna	Image/Video/Audio	Representation-based	Stable Diffusion /AudioLM/Zeroscope
QdP	Cross-modal Understanding	Visual/Audio	Continuous Encoding	CLIP Vit/CLAP	multi-branch	dot attention+Linear	-	injection	DifBERT-V2-XLarge	-	-	-
Uni-Moe	Cross-modal Understanding	Visual/Audio	Continuous Encoding	CLIP Vit/CLAP	multi-branch	Q-former+Linear/MLP	-	concatenate	LLaMA	-	-	-
MoGPT	Cross-modal Understanding	Audio/Motion	Discrete Encoding	/MoGPT's Motion tokenizer	embedding	Extend Vocabulary	-	concatenate	TS	Music/Motion	Modality-Token-based	Jukebox de-tokenizer /MoGPT's Motion de-tokenizer
Emotion-LLaMA	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	CLIP Vit /Juicebox	multi-branch	MLP	-	concatenate	LLaMA-2	-	-	-
EmpathyEar	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	/Juicebox	uni-branch	Linear	-	concatenate	Vicuna	Image/Video/Audio	Representation-based	StyleTTS2/EAT
sales-ALMON	Cross-modal Understanding	Visual/Audio	Continuous Encoding	Vit/Keats	multi-branch	Q-former+Linear	-	concatenate	Vicuna	-	-	-
Merakai	Cross-modal Understanding	Visual/Audio	Continuous Encoding	CLIP Vit/Whisper	multi-branch	MLP	-	concatenate	LLaMA-2	-	-	-
InterOmega	Cross-modal Understanding	Visual/Audio	Continuous Encoding	Intern Vit/Whisper	multi-branch	MLP	-	concatenate	InternLM-2.5	-	-	-
SynceLM	Cross-modal Understanding	Visual/Audio	Hybrid Encoding	SigCLIP Vit /X-Reasons	multi-branch	MLP	Extend Vocabulary	concatenate	OPT	-	-	-
UnifiedMLM	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	/ImageBind-Audio	multi-branch	Q-former+Linear	-	concatenate	OPT	Image/Video/Audio	Text-based	Instruct pic2pix/Affusion/ModelScope GPT-SOTuS
VITA	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	CLIP Vit	multi-branch	MLP	-	concatenate	Mixtral	Speech	Text-based	StyleTTS2/EAT
OcclLMa	Cross-modal Understanding & Generation	3D/Action	Discrete Encoding	OcclLMa-3D tokenizer	embedding	-	Extend Vocabulary	concatenate	LLaMA-3.1	Action/3D	Modality-Token-based	OcclLMa-3D de-tokenizer /OcclLMa-3 Action de-tokenizer
Llama-AVR	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	OcclLMa-3D Action tokenizer	embedding	-	Extend Vocabulary	concatenate	LLaMA-3.1	-	-	-
MIO	Cross-modal Understanding & Generation	Visual/Audio	Discrete Encoding	SEED tokenizer	uni-branch	MLP	-	concatenate	LLaMA-3.1	Yi	Image/Speech	Modality-Token-based /Speech de-tokenizer
EMOVA	Cross-modal Understanding & Generation	Visual/Audio	Hybrid Encoding	Intern Vit	multi-branch	C-Abstractor	Extend Vocabulary	concatenate	LLaMA-3.1	Speech	Modality-Token-based	EMOVA's S2U de-tokenizer MotionRvQ de-tokenizer /Encoder de-tokenizer
LLM Gesticulate	Cross-modal Understanding & Generation	Audio/Motion	Discrete Encoding	MotionRvQ Tokenizer	embedding	-	Extend Vocabulary	concatenate	Qwen-1.5	Audio/Motion	Modality-Token-based	MotionRvQ de-tokenizer /Encoder de-tokenizer
Baichuan-Omni	Cross-modal Understanding	Visual/Audio	Continuous Encoding	SigCLIP Vit/Whisper	multi-branch	CNN+MLP/Conv+MLP	-	concatenate	Vicuna	-	-	-
EAGLE	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	SigCLIP Vit/Whisper	uni-branch	MLP	-	concatenate	LLaMA-2	Image/Video/Audio	Text-based	StableDiffusion /AudioLM/Zeroscope
Spider	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	CLIP Vit/ImageBind-Audio	multi-branch	MLP	-	concatenate	LLaMA-2	Image/Video/Audio	Text-based	StableDiffusion /AudioLM/Zeroscope
Med-2E3	Cross-modal Understanding & Generation	Visual/3D	Continuous Encoding	CLIP Vit/Whisper	multi-branch	Q-former+Linear/MLP	-	concatenate	Phi	-	-	-
LongVAL-E-LLM	Cross-modal Understanding	Visual/Audio	Continuous Encoding	CLIP Vit/Beats/Whisper	multi-branch	MLP	-	concatenate	Vicuna	-	-	-
SOLAMI	Cross-modal Understanding	Audio/Motion	Discrete Encoding	/SOLAMI's MotionTokenizer	embedding	-	Extend Vocabulary	concatenate	Vicuna	Speech/Motion	Modality-Token-based	/SOLAMI's Motion de-tokenizer
MuMu-LLaMA	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	Vit/ViT/MID-MERT	multi-branch	Conv+MLP/Conv+Rnn+MLP	-	injection	LLaMA-2	Music	Representation-based	MusicGen
GMA	Cross-modal Understanding & Generation	Visual/Action	Hybrid Encoding	SigCLIP Vit/GMA Action tokenizer	multi-branch	MLP	Overwrite Vocabulary	concatenate	Qwen-2	Action	Modality-Token-based	GMA's Action de-tokenizer
Lyra	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	DFNCLIP Vit/Whisper	multi-branch	MLP	-	concatenate	Qwen-2	Audio	Representation-based	LLaMA-Omni's audio decoder
VITA-1.5	Cross-modal Understanding & Generation	Visual/Audio	Continuous Encoding	InternVITA's Audio Encoder	multi-branch	MLP/CNN+MLP	-	concatenate	Mixtral	Speech	Representation-based	TiCodec decoder
EmpatheticLM	Cross-modal Understanding	Visual/Audio	Continuous Encoding	CLIP Vit/Whisper	multi-branch	Q-former+Linear	-	concatenate	Qwen-2.5	-	-	-

Table 1: **The architectures of mainstream OmniMLMs.** The architectures of 70 Omni-MLMs are displayed by encoding, alignment, interaction, and generation.

Name	Type	Modality	#Sample
MSCOCO (Lin et al., 2014)	X-Text	Image,Text	620K
Visual Genome (Krishna et al., 2017b)	X-Text	Image,Text	4.5M
Flickr30k (Plummer et al., 2015)	X-Text	Image,Text	158K
SBU (Ordonez et al., 2011)	X-Text	Image,Text	1M
DCI (Urbanek et al., 2024)	X-Text	Image,Text	7.8K
BLIP-Capfilt (Li et al., 2022c)	X-Text	Image,Text	129M
AI Challenger captions (Wu et al., 2017)	X-Text	Image,Text	1.5M
Wukong Captions (Gu et al., 2022)	X-Text	Image,Text	101M
CC12M (Changpinyo et al., 2021)	X-Text	Image,Text	12.4M
CC3M (Sharma et al., 2018)	X-Text	Image,Text	3.3M
LAION-5B (Schuhmann et al., 2022)	X-Text	Image,Text	5.9B
Redcaps (Desai et al., 2021)	X-Text	Image,Text	12M
LAION-COCO (Schuhmann et al., 2022b(b))	X-Text	Image,Text	600M
LAION-CAT (Radenovic et al., 2023)	X-Text	Image,Text	440M
LAION-AESTHETICS (Schuhmann et al., 2022b(a))	X-Text	Image,Text	120M
ShareGPT4V (Chen et al., 2024e)	X-Text	Image,Text	1.2M
LAION-115M (Schuhmann et al., 2021)	X-Text	Image,Text	115M
Journeydb (Sun et al., 2023b)	X-Text	Image,Text	4.4M
Multimodal c4 (Zhu et al., 2023b)	X-Text-X	Image,Text	43.3M
OBELICS (Laurençon et al., 2023)	X-Text-X	Image,Text	141M
Panda-70M (Chen et al., 2024f)	X-Text	Video,Text	70M
Webvid2M (Bain et al., 2021)	X-Text	Video,Text	2M
Valley-Pretrain-703k (Luo et al., 2023a)	X-Text	Video,Text	703K
Webvid10M (Bain et al., 2021)	X-Text	Video,Text	10M
YT-Temporal (Zellers et al., 2022)	X-Text	Video,Text	180M
ActivityNet Captions (Krishna et al., 2017a)	X-Text	Video,Text	100K
InterVid (Wang et al., 2024d)	X-Text	Video,Text	10M
MSRVTT (Xu et al., 2016)	X-Text	Video,Text	200K
ShareGemini (Share, 2024)	X-Text	Video,Text	530K
AudioSet (Gemmeke et al., 2017)	X-Text	Audio,Text	2.1M
Clotho (Drossos et al., 2020)	X-Text	Audio,Text	5k
Auto-ACD (Sun et al., 2024b)	X-Text	Audio,Text	1.5M
AudioCap (Kim et al., 2019)	X-Text	Audio,Text	46k
WavCaps (Mei et al., 2024)	X-Text	Audio,Text	403K
AISHELL-1 (Bu et al., 2017)	X-Text	Audio,Text	128K
AISHELL-2 (Du et al., 2018)	X-Text	Audio,Text	1M
Gigaspeech (Chen et al., 2021)	X-Text	Speech,Text	—
Common Voice (Ardila et al., 2020)	X-Text	Speech,Text	—
MLS (Pratap et al., 2020)	X-Text	Speech,Text	—
Music caption (Zhan et al., 2024)	X-Text	Music,Text	100M
Cap3D (Luo et al., 2023b)	X-Text	3D,Text	1M
Objaverse (Deitke et al., 2023)	X-Text	3D,Text	800K
ScanRefer (Chen et al., 2020a)	X-Text	3D,Text	51.5K
Normal Caption (Han et al., 2024a)	X-Text	Normal,Text	0.5M
Depth Caption (Han et al., 2024a)	X-Text	Depth,Text	0.5M
NSD (Allen et al., 2022)	X-Text	fMRI,Text	9K
Ego4d (Grauman et al., 2022)	X-Text	Video,IMU,Text	528k
PU-VALOR (Tang et al., 2024a)	X-Y-Text	Video,Audio,Text	114K
VALOR (Chen et al., 2023e)	X-Y-Text	Video,Audio,Text	16k
VAST (Chen et al., 2023f)	X-Y-Text	Video,Audio,Text	414k
VIDAL (Zhu et al., 2024a)	X-Y-Text	Video, Thermal, Depth, Audio	10M
TVL (Fu et al., 2024c)	X-Y-Text	Image,Touch,Text	44K
M3D-Cap (Bai et al., 2024a)	X-Y-Text	Image,3D,Text	115K

Table 2: **The statistics for alignment datasets in Omni-MLLMs**, including single non-linguistic modality text pairing data (X-Text), multiple non-linguistic modalities text pairing data (X-Text-Y), and single non-linguistic modality text interleaved data (X-Text-X).

Name	Source	Task	Modality	Construction Method	#Sample
XLLM's SFT (Chen et al., 2023a)	MiniGPT-4, AISHELL-2, VSDial-CN, ActivityNet Caps	Uni-Modal Understanding,Cross-Modal Understanding	Image,Video,Audio,Text	Template Instructionalization, T2X generation	10k
ChatBridge's SFT (Zhao et al., 2023b)	MSRVTT, AudioCaps, VQA2, VG-QA...	Uni-Modal Understanding,Cross-Modal Understanding	Image,Video,Audio,Text	Template Instructionalization, GPT generation	4.4M+209k
Macaw-LLM (Lyu et al., 2023)	MSCOCO, Charades, AVSD, VG-QA...	Uni-Modal Understanding,Cross-Modal Understanding	Image,Video,Audio,Text	GPT generation	69K+50K
BuboGPT's SFT	LLaVA, Clotho, VGGSS	Uni-Modal Understanding,Cross-Modal Understanding	Image,Audio,Text	Template Instructionalization, GPT generation	196K
NextGPT's SFT (Wu et al., 2024b)	WebVid, CC3M, AudioCap, YouTube...	Uni-Modal Understanding,Cross-Modal Understanding, Uni-Modal Generation,Cross-Modal Generation	Image,Video,Audio,Text	Template Instructionalization, GPT generation+retrieval, T2X generation	20K
AnyMal's SFT (Moon et al., 2024)	-	Uni-Modal Understanding	Image,Video,Audio,Text	Manual Annotation, GPT generation	210K
FAVOR's SFT (Sun et al., 2023a)	LLaVA, MSCOCO, Ego4D, LibriSpeech...	Uni-Modal Understanding,Cross-Modal Understanding	Image,Video,Audio,Text	Template Instructionalization, GPT generation	-
LEO's SFT (Huang et al., 2024)	ScanQA, SQAD3, 3RScan, CLIPort...	Uni-Modal Understanding,Cross-Modal Understanding	Image,3D,Text,Action	Template Instructionalization, GPT generation	220k
CoDi-2's SFT (Tang et al., 2024b)	MIMIC-IT, LAION-400M, AudioSet, Webvid...	Uni-Modal Understanding,Uni-Modal Generation	Image,Audio,Text	Template Instructionalization	-
X-InstructBLIP's SFT (Panagopoulou et al., 2024)	MSCOCO, Clotho, MSVD, Cap3D...	Uni-Modal Understanding	Image,Video,Audio,3D,Text	Template Instructionalization, GPT generation	1.6M
OneLLM's SFT (Han et al., 2024a)	LLaVA-150K, Clotho, Ego4D, NSD...	Uni-Modal Understanding	Image,Video,Audio,3D,ImU, Depth,fMRI,Normal,Text	Template Instructionalization, T2X generation	2M
AVLLM's SFT (Shu et al., 2023)	ACAV100M, VGGSound, WebVid2M, WavCaps...	Uni-Modal Understanding,Cross-Modal Understanding	Video,Audio,Text	GPT generation	1.4M
Uni-IO2's SFT (Lu et al., 2024a)	CC3M, AudioSet, Webvid3m, Omni3D...	Uni-Modal Understanding,Cross-Modal Generation	Image,Video,Audio,Text	Template Instructionalization, GPT generation	775m
ModaVerse's SFT (Wang et al., 2024c)	-	Uni-Modal Generation	Image,Video,Audio,Text	GPT generation	2M
REAMO's SFT (Zhang et al., 2024b)	-	Uni-Modal Understanding,Cross-Modal Understanding	Image,Video,Audio,Text	Template Instructionalization	10K
GroundingGPT's SFT (Li et al., 2024h)	Flickr30K, VCR, ActivityNet Captions, Clotho...	Uni-Modal Understanding	Image,Video,Audio,Text	GPT Instructionalization	1M
AnyGPT's SFT (Du et al., 2018)	-	Uni-Modal Understanding,Cross-Modal Understanding	Image,Audio,Text	GPT Instructionalization, T2X generation	208K
CAT's SFT (Ye et al., 2024b)	VGGSound, AVQA, VideoInstruct100K...	Cross-Modal Understanding	Video,Audio,Text	GPT Instructionalization	100K
AVicuna's SFT (Tang et al., 2024a)	UnAV-100, VideoInstruct100K, ActivityNet Captions, DiDeMo	Uni-Modal Understanding,Cross-Modal Understanding	Video,Audio,Text	Template Instructionalization	49K
M3DBench's SFT (Li et al., 2024b)	Scanne, ScanRefer, ShapeNet...	Uni-Modal Understanding,Cross-Modal Understanding	Image,3D,Text	Template Instructionalization, //GPT Instructionalization	320k
Uni-Moe's SFT (Li et al., 2024f)	LLaVA-Instruct-150K, LibriSpeech, Video100K...	Uni-Modal Understanding,Cross-Modal Understanding	Image,Video,Audio,Text	Template Instructionalization, T2X generation	874K
X-VILA's SFT (Ye et al., 2024a)	WebVid, ActivityNetCaption, LLaVA-Instruct-150K...	Uni-Modal Understanding,Cross-Modal Understanding, Uni-Modal Generation,Cross-Modal Generation	Image,Video,Audio,Text	Template Instructionalization	-
EMOVA's SFT (Chen et al., 2024c)	ShareGPT-40, MSCOCO, LLaVA-Instruct-150K...	Uni-Modal Understanding,Cross-Modal Understanding, Uni-Modal Generation,Cross-Modal Generation	Image,Audio,Text	Template Instructionalization, GPT Instructionalization, T2X generation	4.4M
VideoLLaMA2's SFT (Cheng et al., 2024b)	AVQA, AVSD, MusicCaps...	Uni-Modal Understanding,Cross-Modal Understanding	Image,Video,Text	Template Instructionalization	1.5M
PathWeave's SFT (Yu et al., 2024a)	VQA2, MSRVTT, Cap3D...	Uni-Modal Understanding	Image,Video,Audio,3D,Depth	Template Instructionalization, T2X generation	23.2M
Spider's SFT (Lai et al., 2024)	AudioCap, CC3M, Webvid...	Cross-Modal Understanding,Cross-Modal Generation	Image,Video,Audio,Text	Template Instructionalization, GPT Generation	-
GMA's SFT (Szot et al., 2024b)	Meta-World, CALVIN, Manikill...	Uni-Modal Understanding,Cross-Modal Understanding, Cross-Modal Generation	Image,Text,Action	Template Instructionalization	2.2M
OCTAVIUS's SFT (Chen et al., 2024g)	MSCOCO, Bamboo, ScanNet...	Uni-Modal Understanding	Image,3D,Text	Template Instructionalization, GPT Generation	-
Lyra's SFT (Zhong et al., 2024)	Collected YouTube's Audio LibriSpeech, AudioCaps, LLaVA-Instruct-150K...	Uni-Modal Understanding,Cross-Modal Understanding, Uni-Modal Generation	Image,Audio,Text	GPT Generation, T2X Generation	1.5M
video-SALMONN's SFT (Sun et al., 2024a)	VGG-SS, AVSBench, AVQA,MUSIC-AVQA...	Uni-Modal Understanding,Cross-Modal Understanding	Video,Audio,Text	Template Instructionalization, T2X Generation	-
Meerkat's SFT (Chowdhury et al., 2024)	Cross-Modal Understanding	Video,Audio,Text	Template Instructionalization, GPT Generation	3M	
VITA's SFT (Fu et al., 2024b)	ShareGPT4V,LLaVA-Instruct-150K, ShareGPT40, ShareGemini...	Uni-Modal Understanding,Cross-Modal Understanding	Image,Video,Audio,Text	T2X Generation	-
Baichuan-omni's SFT (Li et al., 2024c)	vFLAN, VideoInstruct100K...	Uni-Modal Understanding,Cross-Modal Understanding	Image,Video,Audio,Text	T2X Generation	-
LongVALE-LLM's SFT (Geng et al., 2024)	LongVALE	Uni-Modal Understanding,Cross-Modal Understanding	Video,Audio,Text	GPT Generation	25.4K
UnifiedMLLM's SFT (Li et al., 2024g)	LISA,SmartEdit...	Uni-Modal Understanding,Cross-Modal Understanding, Uni-Modal Generation,Cross-Modal Generation	Image,Video,Audio,Text	Template Instructionalization, GPT Generation	100K
Dolphin's SFT (Anonymous, 2024)	AVQA,Flickr-SoundNet, VGGSound,LLP...	Uni-Modal Understanding,Cross-Modal Understanding	Video,Audio,Text	Template Instructionalization, GPT Generation	-

Table 3: **The statistics for OmniMLLM's Instruction Data**, including the data sources, interaction forms, involved modalities, and construction methods.

Name	Capability Category	Modality	Specific-Task	Metrics
VQA v2 (Goyal et al., 2017)	Unimodal Understanding	Image,Text	QA	Acc
CoQA (Hwang et al., 2019)	Unimodal Understanding	Image,Text	QA	Acc
DocVQA (Mishra et al., 2021)	Unimodal Understanding	Image,Text	QA	Acc
IcoVQA (Lu et al., 2021)	Unimodal Understanding	Image,Text	QA	Acc
OCR-VQA (Mishra et al., 2019)	Unimodal Understanding	Image,Text	QA	Acc
STVQA (Bilen et al., 2019)	Unimodal Understanding	Image,Text	QA	Acc
VSR (Liu et al., 2023a)	Unimodal Understanding	Image,Text	QA	Acc
Harold-Mem (Kumar et al., 2020)	Unimodal Understanding	Image,Text	QA	AUC
OKVQA (Mao et al., 2019)	Unimodal Understanding	Image,Text	QA	Acc
VisWiz (Gurari et al., 2018)	Unimodal Understanding	Image,Text	QA	Acc
TextVQA (Singh et al., 2019)	Unimodal Understanding	Image,Text	QA	Acc
nocap (Agrawal et al., 2019)	Unimodal Understanding	Image,Text	Caption	CIDEr
ScienceQA (Lu et al., 2022)	Unimodal Understanding	Image,Text	QA	Acc
MSCOCO Caption (Lin et al., 2014)	MSCOCO Caption	Image,Text	Caption	CIDEr,BLEU
Flicker30k (Anderson et al., 2015)	Unimodal Understanding	Image,Text	Caption	CIDEr
Visual Dialog (Das et al., 2017)	Unimodal Understanding	Image,Text	Dialogue	MRR
RefCOCO (Yu et al., 2016)	Unimodal Understanding	Image,Text	Grounding	Acc
RefCOCO+ (Yu et al., 2016)	Unimodal Understanding	Image,Text	Grounding	Acc
RefCOCOg (Mao et al., 2016)	Unimodal Understanding	Image,Text	Grounding	Acc
A-vokqa (Schwenk et al., 2022)	Unimodal Understanding	Image,Text	QA	Acc
POPE (Li et al., 2023)	Unimodal Understanding	Image,Text	Hallucination	Acc
ITIFX (Zhou et al., 2012)	Unimodal Understanding	Image,Text	OCR	WAC(word ACC)
IC13 (Karatzas et al., 2013)	Unimodal Understanding	Image,Text	OCR	WAC(word ACC)
IC15 (Karatzas et al., 2015)	Unimodal Understanding	Image,Text	OCR	WAC(word ACC)
Total-Text (Chng and Chan, 2017)	Unimodal Understanding	Image,Text	OCR	WAC(word ACC)
CUTE80 (Risnumawan et al., 2014)	Unimodal Understanding	Image,Text	OCR	WAC(word ACC)
SVTP (Wang et al., 2011)	Unimodal Understanding	Image,Text	OCR	WAC(word ACC)
SVTP (Phu et al., 2013)	Unimodal Understanding	Image,Text	OCR	WAC(word ACC)
COCO-Text (Ved et al., 2016)	Unimodal Understanding	Image,Text	OCR	WAC(word ACC)
MMB (Xu et al., 2016)	Unimodal Understanding	Image,Text	Comprehensive Benchmark	GPT-ACC
MME (Fu et al., 2023)	Unimodal Understanding	Image,Text	Comprehensive Benchmark	GPT ACC
LLVA-Bench (Liu et al., 2023c)	Unimodal Understanding	Image,Text	Comprehensive Benchmark	GPT ACC
Mmme (Yue et al., 2024)	Unimodal Understanding	Image,Text	Comprehensive Benchmark	GPT ACC
SEED (Ge et al., 2023)	Unimodal Understanding	Image,Text	Comprehensive Benchmark	GPT ACC
MM-Vet (Yu et al., 2024d)	Unimodal Understanding	Image,Text	Comprehensive Benchmark	GPT ACC
ActoVQA (Xu et al., 2019)	Unimodal Understanding	Image,Text	Comprehensive Benchmark	GPT ACC
MSRV-TT-QA (Xu et al., 2016)	Unimodal Understanding	Video,Text	QA	Acc
MSVD-QA (Xu et al., 2017)	Unimodal Understanding	Video,Text	QA	Acc
How2QA (Li et al., 2020)	Unimodal Understanding	Video,Text	QA	Acc
NECTQQA (Xiao et al., 2021)	Unimodal Understanding	Video,Text	QA	Acc
STAR (Wu et al., 2021)	Unimodal Understanding	Video,Text	QA	Acc
MSVD-Caption (Xu et al., 2017)	Unimodal Understanding	Video,Text	CIDEr	
VATEX (Cheng et al., 2019)	Unimodal Understanding	Video,Text	Caption	CIDEr
MSRV-TT-Caption (Xu et al., 2016)	Unimodal Understanding	Video,Text	Caption	CIDEr,BLEU
Video-ChaiGPT Benchmark (Maa et al., 2024)	Unimodal Understanding	Video,Text	Comprehensive Benchmark	GPT-ACC,GPT Score
Kinetics-400	Unimodal Understanding	Video,Text	Classification	Acc
Perception test (Patraucean et al., 2023)	Unimodal Understanding	Video,Text	Comprehensive Benchmark	GPT ACC
EgoSchema (Margallam et al., 2023)	Unimodal Understanding	Video,Text	Comprehensive Benchmark	GPT ACC
Mvbench (Li et al., 2024a)	Unimodal Understanding	Video,Text	Comprehensive Benchmark	GPT ACC
VideosMME (Xu et al., 2024a)	Unimodal Understanding	Video,Text	Comprehensive Benchmark	GPT ACC
Charon-STA (Song et al., 2016)	Unimodal Understanding	Video,Text	Grounding	IoU
AudioCaps (Kim et al., 2019)	Unimodal Understanding	Audio,Text	Caption	CIDEr,SPICE,METEOR,BLEU,SPIDER
CllothoQA (Liping et al., 2022)	Unimodal Understanding	Audio,Text	QA	Acc
VocalSound (Gong et al., 2022)	Unimodal Understanding	Audio,Text	QA	Acc
Cllotho v1 (Drossos et al., 2020)	Unimodal Understanding	Audio,Text	Caption	CIDEr
Cllotho v2 (Drossos et al., 2020)	Unimodal Understanding	Audio,Text	Caption	CIDEr
ESCC50 (Pezak, 2015)	Unimodal Understanding	Audio,Text	Classification	Acc
LibriSpeech (Hannun et al., 2015)	Unimodal Understanding	Audio,Text	ASR	WER
AIHELL-2 (Dax et al., 2018)	Unimodal Understanding	Audio,Text	ASR	WER
WenoteSpeech (Zhang et al., 2022a)	Unimodal Understanding	Audio,Text	ASR	WER
MusicCap (Agostinelli et al., 2023)	Unimodal Understanding	Audio,Text	Caption	CLAP Score
TUT2017 (Mesaros et al., 2016)	Unimodal Understanding	Audio,Text	Classification	Acc
EHSI (Li et al., 2024d)	Unimodal Understanding	Audio,Text	QA	Acc
Cap3D-Caption (Deitke et al., 2023b)	Unimodal Understanding	3D,Text	Caption	CIDEr
Objavise-Caption (Deitke et al., 2023)	Unimodal Understanding	3D,Text	Caption	METEOR,ROUGE,BLEU
Cap3D-QA (Luo et al., 2023b)	Unimodal Understanding	3D,Text	QA	Acc
Obiverse Classification (Deitke et al., 2023)	Unimodal Understanding	3D,Text	Classification	GPT ACC
Modelnet40 (Wu et al., 2015)	Unimodal Understanding	3D,Text	Classification	Acc
ScanRefer (Chen et al., 2020a)	Unimodal Understanding	3D,Text	Classification	Acc
Nr3D (Achlioptas et al., 2020)	Unimodal Understanding	3D,Text	Classification	Acc
SQARL (Xu et al., 2020)	Unimodal Understanding	3D,Text	Classification	Acc
Sun3DQA (Xu et al., 2022)	Unimodal Understanding	Depth,Text	Classification	Acc
SUN-RGB-D (Song et al., 2015)	Unimodal Understanding	Depth,Text	Classification	Acc
NYUV2 (Silberman et al., 2012)	Unimodal Understanding	Depth,Text	Normal,Text	Classification
SUN-RGB-D_generated Nomral (Han et al., 2024a)	Unimodal Understanding	Normal,Text	Classification	Acc
NYUV2_generated Nomral (Han et al., 2024a)	Unimodal Understanding	Normal,Text	Classification	Acc
ThermalQQA (Yu et al., 2024c)	Unimodal Understanding	Thermal,Text	QA	Acc
DocQQA (Yu et al., 2024c)	Unimodal Understanding	Touch,Text	QA	Acc
Ego4D (Garg et al., 2022)	Unimodal Understanding	IMU,Text	Caption	CIDEr,ROUGE
NSD (Allen et al., 2022)	Unimodal Understanding	IMU,Depth,Map,Text	Caption	CIDEr,ROUGE
MSCOCO (Lin et al., 2014)	Unimodal Generation	Image,Text	TX2X Edit	FID,CLIPSIM
MSRVTT (Xu et al., 2016)	Unimodal Generation	Video,Text	T2X Generate	CLIPSIM
AudioCaps (Kim et al., 2019)	Unimodal Generation	Audio,Text	T2X Generate,TX2XX Edit	FAD
DAVIS (Perazzi et al., 2016)	Unimodal Generation	Video,Text	T2X Edit	CLIPSIM
UCF-101 (Soomro et al., 2012)	Unimodal Generation	Video,Text	T2X Generate	FID,CLIPSIM
Event-50 (Xu et al., 2024c)	Unimodal Generation	Video,Text	T2X Generate,TX2XX Edit	WER,MCD
VCTK (Vassau et al., 2017)	Unimodal Generation	Audio,Text	T2X Generate	FAD
MusicCap (Agostinelli et al., 2023)	Unimodal Generation	Audio,Text	T2X Generate	CLIP-LCLIP-T,DINO
Dreambench (Ruiz et al., 2023)	Crossmodal Understanding	Image,Text	QA	Acc
MUSIC-AVQA (Li et al., 2022b)	Crossmodal Understanding	Video,Audio,Text	Dialogue	CIDEr,BLEU
AVSD (Al-Antri et al., 2018)	Crossmodal Understanding	Image,Video,Audio,Text	Comprehensive Benchmark	Acc
RACE (Liu et al., 2024)	Crossmodal Understanding	Image,Video,Audio,Text	Comprehensive Benchmark	Acc
VALOR-Caption (Chen et al., 2023e)	Crossmodal Understanding	Video,Audio,Text	Caption	CIDEr,BLEU
MMBench-Audio (Li et al., 2024f)	Crossmodal Understanding	Image,Audio,Text	Comprehensive Benchmark	Acc
AVQA (Li et al., 2022a)	Crossmodal Understanding	Video,Audio,Text	QA	Acc
MCUB (Chen et al., 2024a)	Crossmodal Understanding	Image,Video,Audio,3D,Text	Comprehensive Benchmark	Acc
DisCRn (Panagopoulou et al., 2024)	Crossmodal Understanding	Image,Video,Audio,3D	Comprehensive Benchmark	Acc
OmniXR (Chen et al., 2024d)	Crossmodal Understanding	Image,Video,Audio,Text	Comprehensive Benchmark	Acc
Curve (Leng et al., 2024)	Crossmodal Understanding	Image,Video,Audio,Text	Comprehensive Benchmark	Acc
ISOQA (Xu et al., 2024)	Crossmodal Understanding	Image,Video,Audio,Text	Hallucination	Acc
VGGSound (Chen et al., 2020b)	Crossmodal Understanding	Video,Audio,Text	QA	Acc
VATEX (Wang et al., 2019)	Crossmodal Understanding	Video,Audio,Text	Caption	CIDEr
UnAV-100 (Geng et al., 2023)	Crossmodal Understanding	Video,Audio,Text	Ground	IoU
LJP (Tian et al., 2020)	Crossmodal Understanding	Video,Audio,Text	Ground	IoU
PretzelCaption-QA (Sun et al., 2024a)	Crossmodal Understanding	Video,Audio,Text	QA	Acc
Long-VC (Chen et al., 2024)	Crossmodal Understanding	Video,Audio,Text	Caption	CIDEr
TVL Benchmark (Fu et al., 2024c)	Cross-modal Understanding	Touch,Image,Text	QA	Acc
AVED (Sun et al., 2023a)	Cross-modal Understanding	Image,Video,Audio,Text	Comprehensive Benchmark	ACC,METEOR,SPIDER,WER
XioX Benchmark (Ye et al., 2024a)	Crossmodal Understanding	Image,Video,Audio,Text	Comprehensive Benchmark	X-to-X Alignment Score

Table 4: An overview of benchmarks and tasks of Omni-MLLMs, including the abilities being evaluated, the involved modalities, specific tasks, and evaluation metrics.

Model	LLM	Uni-Modal Understanding								Uni-Modal Generation			Cross-Modal Understanding			
		MSVD-QA	MSRVTT-QA	VQA ^{v2 test}	Flickr	MMB ^{cen}	AudioCaps ^{cap test}	ClothoAQA	Objaverse	COCO ^{sem}	AudioCaps ^{gen}	MSRVTT ^{sem}	VGGSS	AVSD	MUSIC-AVQA	AVQA
Omni-MLLMs																
eP-ALM	OPT-2.7B	38.4	38.51	54.47	—	—	61.86	—	—	—	—	—	—	—	—	
ChatBridge 13B	Vicuna-13B	45.3	—	—	82.5	—	—	—	—	—	—	—	—	43	—	
PandaGPT	Vicuna-13B	46.7	23.7	—	—	—	—	—	—	—	—	32.7	26.1	33.7	79.8	
Video-LLaMA	Vicuna-7B	51.6	29.6	—	—	—	—	—	—	—	—	40.8	36.7	36.6	81	
Macaw-LLM	LLaMA-7B	42.1	25.5	—	—	3.84	33.3	—	—	—	—	36.1	34.3	31.8	78.7	
ImageBind-LLM	LLaMA-7B	—	—	23.49	—	—	10.3	31	—	—	—	—	—	39.72	54.26	
NExt-GPT	Vicuna-7B	64.5	58.4	66.7	84.5	58	81.3	—	—	10.07	8.67	31.97	—	—	—	
AnyMAL 13B	LLaMA2-13B	—	—	59.6	—	—	—	—	—	—	—	—	—	—	—	
AnyMAL 70B	LLaMA2-70B	—	—	64.2	95.9	—	77.8	—	—	—	—	—	—	—	—	
X-InstructBLIP 7B	Vicuna-7B	51.7	41.3	30.61	82.1	8.96	67.9	15.4	50	—	—	—	—	28.1	—	
X-InstructBLIP 13B	Vicuna-13B	49.2	—	—	74.7	—	53.7	21.7	—	—	—	20.3	52.1	44.5	44.23	
OneLLM 7B	LLaMA2-7B	56.5	—	71.6	78.6	60	—	57.9	44.5	—	—	—	—	47.6	—	
AV-LLM	Vicuna-7B	67.3	53.7	—	—	—	35.5	—	—	—	—	47.6	52.6	45.2	—	
UIO-2x1 6.8B	—	52.2	41.5	79.4	—	71.5	48.9	—	—	13.39	5.89	—	—	—	—	
ModaVerse	Vicuna-7B	—	56.5	—	—	—	79.2	—	—	11.24	8.22	30.14	—	—	—	
CREMA 7B	Mistral-7B	—	—	—	—	—	—	—	—	—	—	—	—	52.6	—	
GroundingGPT	Vicuna-7B	67.8	51.6	78.7	—	63.8	—	—	—	—	—	—	—	—	—	
NaiveMC	Vicuna-7B	—	—	—	—	—	—	—	55	—	—	—	—	53.63	80.7	
DAMC	Vicuna-7B	—	—	—	—	—	—	—	60.5	—	—	—	—	57.32	81.31	
AnyGPT	LLaMA2-7B	—	—	—	—	—	—	—	—	—	—	—	—	—	—	
CAT	LLaMA2-7B	—	62.7	—	—	—	—	—	—	—	—	—	48.6	92	—	
AVicuna	Vicuna-7B	70.2	59.7	—	—	—	—	—	—	—	—	—	53.1	49.6	—	
Uni-MoE	LLaMA-7B	55.6	—	66.2	—	69.82	—	32.6	—	—	—	—	—	—	—	
X-VILA 7B	Vicuna-7B	—	—	72.9	—	—	—	—	—	—	—	—	—	—	—	
VideoLLaMA2-7B	Mistral-7B	71.7	—	—	—	—	—	—	—	—	—	—	71.4	57.2	80.9	
Meerkat	Llama-2-7B-Chat	—	—	—	—	—	—	—	—	—	—	—	—	—	87.14	
InternOmni	InternLM-2-Chat-7B	—	—	—	—	81.7	—	—	—	—	—	—	—	—	—	
UnifiedMLLM	Vicuna-7B	—	—	—	—	—	—	—	—	—	—	—	—	—	—	
VITA	Mixtral-8x7B	—	—	—	—	—	71.8	—	—	—	—	—	—	—	—	
EMOVA	LLaMA-3.1-8B	—	—	—	—	—	82.8	—	—	—	—	—	—	—	—	
BaiChuan-omni-7B	—	72.2	—	—	—	76.2	—	—	—	—	—	—	—	—	—	
OMCAT	Vicuna-7B	—	—	—	—	—	—	—	—	—	—	—	49.4	73.8	90.2	
PathWeave-7B	Vicuna-7B	47.8	37.4	—	—	—	64	33.5	—	—	—	—	—	—	—	
Spider	Llama-2-7B	—	—	—	—	—	81.7	—	—	11.23	8.18	30.97	—	—	—	
Specific-MLLMs																
VILA-7B	LLaMA-2-7B	—	—	79.9	74.7	68.9	—	—	—	—	—	—	—	—	—	
VideoChat2	Vicuna-7B	70	54.1	—	—	—	—	—	—	—	—	—	—	—	—	
Qwen-Audio	Qwen-7B	—	—	—	—	—	—	57.9	—	—	—	—	—	—	—	
PointLLM	Vicuna-7B	—	—	—	—	—	—	—	47.5	—	—	—	—	—	—	
Emu-13B	—	—	—	52	—	—	—	—	—	11.66	—	—	—	—	—	
Video-LaVIT	Llama2-7B	73.2	—	80.3	—	67.3	—	—	—	—	30.12	—	—	—	—	

Table 5: **The performance of Omni-MLLMs on different benchmarks.** The selected uni-modal understanding benchmarks include Video-Text2Text (Xu et al., 2017, 2016), Image-Text2Text (Goyal et al., 2017; Plummer et al., 2015; Liu et al., 2024d), Audio-Text2Text (Kim et al., 2019; Lipping et al., 2022), and 3D-Text2Text (Deitke et al., 2023). The chosen uni-modal generation benchmarks include Text2Image (Lin et al., 2014), Text2Video (Xu et al., 2016), and Text2Audio (Kim et al., 2019). The selected cross-modal understanding benchmarks are Image-Audio-Text2Text (Chen et al., 2020b; AlAmri et al., 2018) and Video-Audio-Text2Text (Li et al., 2022b,a).