

Model Composition for Multimodal Large Language Models

Chi Chen^{*1}, Yiyang Du^{*1}, Zheng Fang¹, Ziyue Wang¹, Fuwen Luo¹,
Peng Li^{✉2,4}, Ming Yan³, Ji Zhang³, Fei Huang³, Maosong Sun¹, Yang Liu^{✉1,2,4,5}

¹Dept. of Comp. Sci. & Tech., Institute for AI, Tsinghua University, Beijing, China

²Institute for AI Industry Research (AIR), Tsinghua University, Beijing, China

³Institute of Intelligent Computing, Alibaba Group

⁴Shanghai Artificial Intelligence Laboratory, Shanghai, China

⁵Jiangsu Collaborative Innovation Center for Language Competence, Jiangsu, China

Abstract

Recent developments in Multimodal Large Language Models (MLLMs) have shown rapid progress, moving towards the goal of creating versatile MLLMs that understand inputs from various modalities. However, existing methods typically rely on joint training with paired multimodal instruction data, which is resource-intensive and challenging to extend to new modalities. In this paper, we propose a new paradigm through the model composition of existing MLLMs to create a new model that retains the modal understanding capabilities of each original model. Our basic implementation, NaiveMC, demonstrates the effectiveness of this paradigm by reusing modality encoders and merging LLM parameters. Furthermore, we introduce DAMC to address parameter interference and mismatch issues during the merging process, thereby enhancing the model performance. To facilitate research in this area, we propose MCUB, a benchmark for assessing ability of MLLMs to understand inputs from diverse modalities. Experiments on this benchmark and four other multimodal understanding tasks show significant improvements over baselines, proving that model composition can create a versatile model capable of processing inputs from multiple modalities.¹

1 Introduction

Recent advancements in Multimodal Large Language Models (MLLMs) have established them as the forefront of multimodal learning paradigms (Liu et al., 2023b,a). The prevalent approach involves aligning modality encoders with large language models (LLMs) through extensive modality-text paired data and then fine-tuning with modality-specific instruction data. This paradigm

^{*}Equal contribution.

[✉]Corresponding authors: Peng Li (lipeng@air.tsinghua.edu.cn) and Yang Liu (liuyang2011@tsinghua.edu.cn).

¹Code is available at <https://github.com/THUNLP-MT/ModelCompose>

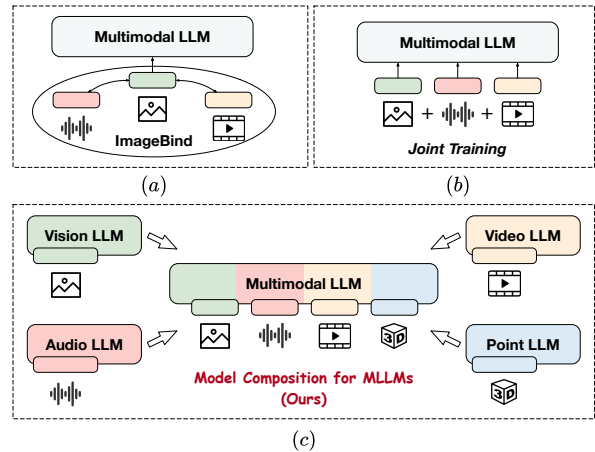


Figure 1: Illustration of various approaches for multimodal large language models: (a) aligning LLM with a multimodal encoder and (b) joint training with multiple modal encoders and (c) our proposed model composition method that creates a versatile model from existing MLLMs through a training-free and extensible process.

has been successfully applied to a wide range of modalities such as image (Liu et al., 2023b; Dai et al., 2023), audio (Gong et al., 2023; Deshmukh et al., 2023; Tang et al., 2024), video (Zhang et al., 2023; Lin et al., 2023), and point cloud (Hong et al., 2023; Xu et al., 2023), resulting in the emergence of a diverse array of MLLMs with unique modal capabilities.

There are also some efforts to enable a single MLLM to handle multiple modalities. One way to achieve this is to align the LLM with a multimodal encoder such as ImageBind (Girdhar et al., 2023) with only image-text data (Figure 1(a)) (Su et al., 2023; Han et al., 2023). This approach leverages the inherent alignment of different modalities within the multimodal encoder, allowing the MLLM to comprehend various modalities to a certain degree. However, the absence of modality-specific instruction data often results in suboptimal performance. Another approach entails the concurrent training of the MLLM with multiple

modality encoders (Figure 1(b)). For example, ChatBridge (Zhao et al., 2023) connects image, video and audio encoders with the LLM through a joint training process with multimodal instruction data (i.e., video-audio chats). This kind of methods show potential but faces two major challenges. First, it is resource-heavy to collect paired data across multiple modalities. Second, adapting these models to new modalities requires additional training, adding to the complexity and resource demands of the development process.

Given the limitations of current approaches, we propose and study a more practical setting: *model composition* for MLLMs (Figure 1(c)). Our primary research question is simple yet fundamental: *Can we create a new MLLM by combining existing MLLMs to inherit their understanding of different modalities without training?* Model composition for MLLMs is advantageous for two key reasons: (1) it eliminates the need for the resource-heavy process of training and gathering multimodal data, and (2) it promises enhanced adaptability, facilitating seamless incorporation of new modalities.

Some recent studies, such as X-InstructBLIP (Panagopoulou et al., 2023), serve as pioneering efforts in model composition for MLLMs. These works primarily train projectors to align different encoders with a single LLM and demonstrate the ability to process multiple modalities concurrently. However, a critical limitation is their applicability only to MLLMs with **frozen** language model weights. This constraint restricts the range of models that can be utilized, and impairs the overall performance of the MLLMs (Zeng et al., 2023).

In this paper, we first propose a framework for model composition for MLLMs. Our implementation, named NaiveMC, is elegantly simple: for the MLLMs to be composed, we directly reuse their modality-specific encoders and merge their LLM parameters. We demonstrate that MLLMs, as long as initialized from the same LLM, can achieve zero-shot multi-modality expansion through this model composition framework, regardless of whether the parameters of the LLM have been fine-tuned.

Furthermore, to mitigate parameter interference in the composition process and optimize the performance of the composite model, we propose DAMC, an advanced framework with parameter **D**ecoupling and **A**djustment for **M**odel **C**omposition. By separating modality-specific parameters from language model parameters during initial MLLM training, DAMC allows for the selec-

tive merging of textual parameters, reducing cross-modal interference. Moreover, DAMC introduces an adaptive parameter adjustment mechanism to ensure optimal compatibility and effectiveness of the composite model, achieving a balanced and efficient multi-modality expansion.

To assess the efficacy of our proposed frameworks, we conduct comprehensive experiments on tasks that require an integrated understanding of inputs from diverse combinations of four prevalent modalities: image, audio, video, and point cloud. To facilitate the research in model composition for MLLMs, we also build MCUB, a benchmark specifically designed for evaluating the capability to concurrently comprehend multiple modalities by identifying commonalities across inputs from various modalities. Experimental results indicate that our frameworks enable the composition of existing MLLMs from different modalities without requiring further training, yielding a versatile and high-performing multimodal model adept at handling any combination of these modalities.

Our contributions are three-fold:

- We propose the concept of model composition for MLLMs, realized through the NaiveMC framework, which allows for seamless integration of different MLLMs without additional training, enabling zero-shot multi-modality expansion.
- We introduce DAMC, an advanced model composition framework that employs parameter decoupling and adaptive adjustment to mitigate parameter interference and optimize composite model performance across multiple modalities.
- We create MCUB, a benchmark designed to evaluate the unified understanding of diverse modalities, and demonstrate the efficacy of our model composition frameworks via extensive experiments on various multimodal understanding tasks and MCUB.

2 Related Work

2.1 Multimodal Large Language Models

Recent advancements have integrated Large Language Models (LLMs) with multiple modalities such as image, video, audio, and point cloud. Approaches like X-LLM (Chen et al., 2023) utilize modality-specific adapters, while ChatBridge (Zhao et al., 2023) employs a Perceiver for each modality. Macaw-LLM (Lyu et al., 2023) uses a

unified alignment module, and NeXT-GPT (Wu et al., 2023) relies on linear projection for multimodal integration. However, these methods often necessitate joint multimodal dataset training, posing scalability and extension challenges. While PandaGPT (Su et al., 2023) and ImageBind-LLM (Han et al., 2023) bypass joint training with ImageBind (Girdhar et al., 2023) as a unified encoder, their performance is limited by a lack of modality-specific instruction data. In contrast, our approach benefits from instruction-following training for each MLLM and integrates their capabilities through a training-free process.

2.2 Model Composition

Model composition integrates different models’ capabilities, primarily through weight space merging strategies (Ilharco et al., 2022; Matena and Raffel, 2022; Jin et al., 2023; Huang et al., 2023; Yu et al., 2024; Wortsman et al., 2022), enhancing fine-tuned models from the same initialization. However, existing research has not fully explored model composition for multimodal LLMs. Prior works on multimodal tasks use linear interpolation for tasks like speech recognition (Sundar et al., 2023). Another work converts model processing images and text separately into one using same parameters for simultaneous processing (Sung et al., 2023). But they do not expand the processing modalities of the model. X-InstructBlip (Panagopoulou et al., 2023) incorporates new modalities by training separate projections for modal encoders and LLM, but its performance is limited by the frozen LLM. Our work advances multimodal capabilities by enabling LLM training and applying model composition techniques to LLM weights, significantly improving the effect of multimodal integration.

3 Methodology

3.1 Task Definition

In this work, we examine the composition of a set of Multimodal Large Language Models (MLLMs), denoted as $\{M_1, M_2, \dots, M_n\}$, each capable of responding to textual queries with different modalities $\mathbf{m} = \{m_1, m_2, \dots, m_n\}$. The core objective is to develop a composition method C that effectively integrates these MLLMs into a singular, more versatile model $M_{\text{compose}} = C(M_1, M_2, \dots, M_n)$. This integration aims to enable the model to process and understand inputs from any combination of modalities in \mathbf{m} . For example, by integrating two

specialized MLLMs—a vision LLM and an audio LLM—the resulting composite model should not only preserve the individual proficiencies of these models in processing images and audio, respectively, but should also acquire a zero-shot capacity for handling inputs that encompass both visual and auditory information simultaneously.

Algorithm 1 Model Composition for MLLMs

Input: Set of MLLMs $\{M_1, M_2, \dots, M_n\}$, each with a set of parameters Θ^i

Output: Integrated parameters Θ_{compose}

```

1: Initialize  $\Theta_{\text{common}}$  and  $\Theta_{\text{compose}}$  as empty sets
2: Define a mapping  $f : \theta \mapsto G$  that maps the
   parameter  $\theta$  to group  $G$  based on functionality
3: for  $i = 1$  to  $n$  do
4:   for each parameter  $\theta$  in  $\Theta^i$  do
5:     if  $\forall j \neq i, \theta \notin \Theta^j$  then
6:       Add  $\theta$  to  $\Theta_{\text{compose}}$ 
7:     else
8:       Assign  $\theta$  to group  $f(\theta)$  in  $\Theta_{\text{common}}$ 
9:     end if
10:  end for
11: end for
12: for each group  $G$  in  $\Theta_{\text{common}}$  do
13:    $\theta_{\text{merge}} = \text{average}(\text{parameters in } G)$ 
14:   Add  $\theta_{\text{merge}}$  to  $\Theta_{\text{compose}}$ 
15: end for
16: return  $\Theta_{\text{compose}}$ 

```

3.2 A Model Composition Framework

When composing models, two critical elements must be taken into account: the **components** and their **weights**. The prevalent MLLMs typically feature two component types: (1) *modal-specific* components, like modal encoders and projectors, which adapt modality inputs to the language embedding space and (2) *modal-agnostic* components that exist in each MLLM, primarily the underlying LLM itself. In our composition framework, we retain all modal-specific components (and their weights) from different MLLMs to handle respective modal inputs, and connect them to the same LLM. In cases where the LLMs have not been adapted during the training of MLLMs (as illustrated in Figure 2(a)), we employ the pre-trained weights of the LLM directly. Conversely, if the LLMs have undergone adaptation in the MLLM training process (Figure 2(b)), we simply average their weights. We name this composition framework NaiveMC and provide a formal procedure of it in Algorithm 1.

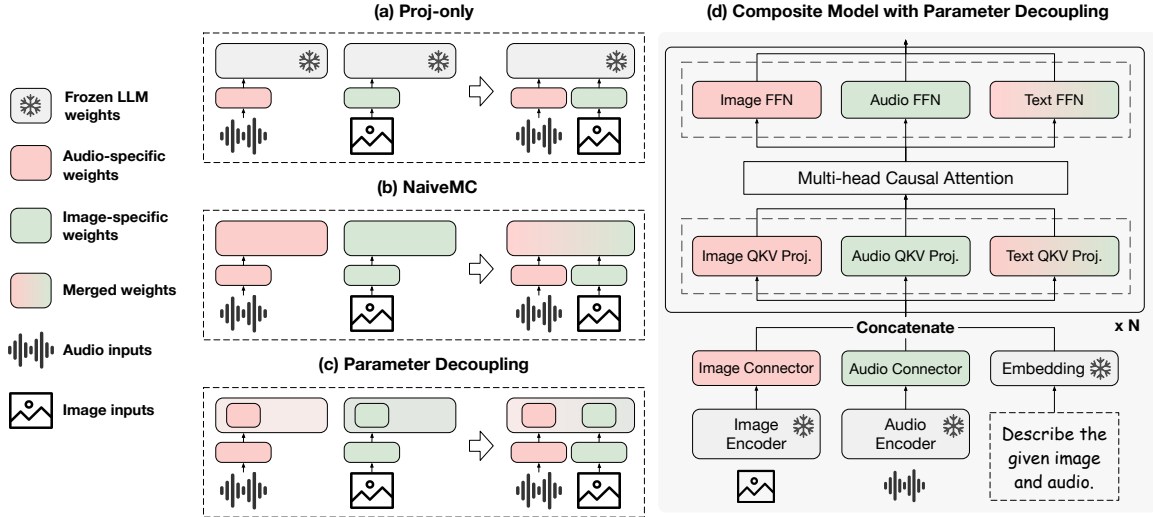


Figure 2: Illustration of the model composition processes with only image and audio modalities are considered for simplicity. (a) and (b) show a basic model composition framework as described in Section 3.2, while (c) and (d) demonstrate model composition with parameter decoupling, as detailed in Section 3.3.

3.3 Parameter Decoupling

In the previously described framework, there arises an unaddressed issue, specifically the potential for parameter interference when merging fine-tuned LLM parameters. During the training of MLLMs, the LLM parameters are optimized for specific modal inputs. This specialized optimization could lead to variations in the LLM parameters of different MLLMs. When these parameters are merged, conflicts may arise, potentially impacting the ability of the model to understand modal inputs.

To alleviate the issue of parameter interference during model composition, we advocate for initially training the MLLMs with a parameter decoupling strategy in the first place. As shown in Figure 2(c) and Figure 2(d), the main idea is to separate the modality processing parameters from those of the language model within MLLMs. For example, in an MLLM M that processes modality m , the input for each attention layer is denoted as $X = [X_m, X_t]$ where X_m and X_t represent the modality-specific and text sequences, respectively. The attention components are computed as follows:

$$Q = [X_m W_m^Q, X_t W_t^Q] \quad (1)$$

$$K = [X_m W_m^K, X_t W_t^K] \quad (2)$$

$$V = [X_m W_m^V, X_t W_t^V] \quad (3)$$

where Q , K , and V represent the queries, keys, and values in the attention mechanism, respectively. Note that two sets of weights $W_m^{(\cdot)}$ and $W_t^{(\cdot)}$ are employed, tailored to the modality-specific and tex-

tual inputs. The attention operation is then applied:

$$[X_m^O, X_t^O] = \text{split}(\text{Attention}(Q, K, V)) \quad (4)$$

which splits the attention output into modality-specific and textual outputs. The final unified output representation X_o is obtained by:

$$X_o = [X_m^O W_m^O, X_t^O W_t^O] \quad (5)$$

ensuring separate processing streams for each modality within the model. Similarly, given input X for each feed-forward layer, the output is:

$$\text{FFN}(X) = [\text{FFN}_m(X_m), \text{FFN}_t(X_t)] \quad (6)$$

where FFN_m and FFN_t are feed-forward layers for modality-specific and text inputs, respectively.

When composing MLLMs that are trained through parameter decoupling, we merge only the text-related parameters, maintaining distinct modality-specific parameters as depicted in Figure 2(d). The composite model functions as a natural extension of the methodology previously described, guaranteeing that inputs from each modality are processed independently with their respective parameters. By doing so, it effectively mitigates the risk of interference from other modalities, ensuring that the composite model maintains high fidelity in processing multimodal data. Please note that after the MLLMs are trained, the composition phase remains emphatically *training-free*.

3.4 Adaptive Parameter Adjustment

Due to variations in data quality and training strategies among different MLLMs, their performance can significantly differ. Thus, when composing these models, employing a simplistic averaging strategy often falls short of achieving optimal results. To enhance model effectiveness, it is crucial to implement adaptive adjustments to the model parameters during the composition process, allowing for better compatibility and flexibility among the different models. Specifically, for N distinct MLLMs $\{M_1, M_2, \dots, M_N\}$ and any parameter θ_i common across these models, where $i = 1 \dots N$, the merged parameter is defined as:

$$\theta_{\text{merge}} = \sum_{i=1}^N \lambda_i \theta_i \quad (7)$$

where λ_i represents the adjustment coefficient. For simplicity, we adopt a uniform adjustment coefficient λ_i for all θ_i in M_i . For models trained using parameter decoupling, we can additionally adjust their modality-specific parameters if needed. The values of these coefficients can be determined with a validation set from target tasks requiring various modal inputs. If such a validation set is not available, a practical alternative is to select the coefficients based on general performance of the model on tasks of each modality. We refer to the updated model composition framework with parameter decoupling and adjustment as DAMC.

3.5 Multimodal Commonality Understanding Benchmark

As previously discussed, DAMC inherently possesses the advantage of scaling up the number of modalities. However, there are few benchmarks available to evaluate the performance across various modalities. To demonstrate the effectiveness of our approach on tasks involving numerous modalities, inspired by Panagopoulou et al. (2023), we introduce a new benchmark called the Multimodal Commonality Understanding Benchmark (MCUB). We provide an example of MCUB in Figure 3. The task of MCUB is to measure the ability of the model to identify commonalities among input entities from diverse modalities and select the most appropriate answer from four given candidates.

We leverage GPT-4 (Achiam et al., 2023) to create MCUB from existing captioning datasets across various modalities. To be specific, we begin by randomly selecting groups of inputs from each modal-

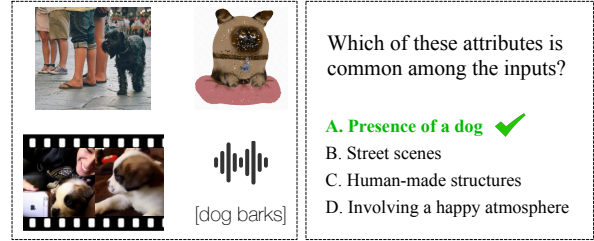


Figure 3: An example of MCUB-4, where the objective is to identify common attributes from inputs across four different modalities.

ity, and retain the groups with the highest semantic similarity. This similarity is quantified by averaging the similarities of their respective captions. Subsequently, GPT-4 is prompted with in-context examples to generate questions, options, and correct answers for each group. In this study, we develop two variants of MCUB: MCUB-3 and MCUB-4. MCUB-4 comprises data entries that include inputs from four modalities—image, video, audio, and point cloud. In contrast, MCUB-3 consists of four subsets, each representing combinations of inputs from any three of these modalities. For more details, please refer to Appendix C.

4 Experiments

4.1 Implementation

In our experimental setup, we explore four modalities: *image*, *audio*, *video* and *point cloud*. To ensure a comprehensive and comparative analysis, we reimplement MLLMs for each modality following previous works (Liu et al., 2023a; Lin et al., 2023; Panagopoulou et al., 2023; Xu et al., 2023). For each of these modalities, we train three versions of MLLMs with same data and hyperparameters but vary the trainability of LLM parameters: a model with a frozen LLM, a fully trainable LLM, and a trainable LLM with parameter decoupling, as illustrated in Figure 2. We employ the LoRA (Hu et al., 2021) technique for efficient LLM training. All models are based on Vicuna-7B-v1.5 (Zheng et al., 2023). More details for training the MLLMs are in Appendix A. We apply parameter adjustment by conducting a search for the optimal λ_i values within the set $[1/N, 2/N, \dots, N/N]$ for composition of N modalities based on the performance of the validation set for each corresponding task. The results for the best-performing parameters are documented in Appendix B.

Task	Method	V	V + I	V + A	V + I + A
MUSIC-AVQA	ChatBridge-13B	-	-	43.00	-
	OneLLM-7B	-	-	47.60	-
	ImageBind-LLM	37.24	38.76	39.72	38.16
	X-InstructBLIP	45.83	41.23	48.34	47.39
	Proj-only	44.93	46.64	46.17	50.21
	NaiveMC	49.00	52.52	50.66	53.63
	DAMC	49.09	53.08	50.91	57.32
AVQA	ImageBind-LLM	51.77	51.65	55.00	54.26
	X-InstructBLIP	41.91	40.42	44.29	44.23
	Proj-only	67.99	66.65	67.65	66.85
	NaiveMC	79.37	79.74	79.82	80.70
	DAMC	79.15	80.30	80.40	81.31

Table 1: Experimental results on zero-shot audio-visual question answering tasks with different combinations of video (V), image (I) and audio (A) inputs. Methods developed in this study are distinguished with a grey background for clarity.

4.2 Evaluation

Datasets and Benchmarks. Our aim is to evaluate the performance of the composite MLLM to process inputs from multiple modalities. In addition to the MCUB benchmark described in Section 3.5, we also include datasets from two key areas for evaluation: (1) *audio-visual question answering* including MUSIC-AVQA (Li et al., 2022) and AVQA (Yang et al., 2022) where image, video and audio inputs are available and (2) *3D object classification* on ModelNet40 (Wu et al., 2015) and Objaverse (Deitke et al., 2023) where image and point cloud inputs are considered.

Baselines. Our evaluation strategy involves comparison with baselines primarily focusing on models that do not leverage training data from multiple modalities. We consider two representative baselines for this purpose: **ImageBind-LLM** that aligns a multimodal encoder, ImageBind, to the LLM with image-text paired data and **X-InstructBLIP** that aligns individual modal encoders to a frozen LLM using instruction data specific to each modality. Regarding model composition, we further compare our DAMC approach against two alternative intermediate baselines as discussed before: **Proj-only** where the LLM is frozen and only the modal encoders and connectors are composed and **NaiveMC** that directly composes MLLMs without adopting our proposed parameter decoupling and adjustment techniques.

Task	Method	Type-I	Type-C
Objaverse	3D-LLM	49.00	41.50
	ImageBind-LLM	31.00	26.50
	X-InstructBLIP	50.00	31.50
	Proj-only	48.00	42.50
	NaiveMC	55.00	59.50
	DAMC	60.50	62.00
ModelNet40	ImageBind-LLM	42.71	42.46
	X-InstructBLIP	61.43	61.14
	Proj-only	62.88	61.99
	NaiveMC	66.00	64.59
	DAMC	70.02	65.24

Table 2: Experimental results on zero-shot 3D object classification tasks using combined point and image inputs (P + I). Following Xu et al. (2023), two different types of prompts are considered: Type-I (“What is this?”) and Type-C (“This is an object of”).

4.3 Main Results

The main experimental results on audio-visual question answering, 3D object classification and MCUB are presented in Table 1, Table 2 and Table 3, respectively. In the tables, the notation “X+Y” represents the combination of different types of modal inputs used in evaluating the models. For example, “V + A” refers to the combination of video (V) and audio (A) inputs being used together as part of the input to the model. For model composition, we compose MLLMs based on the specific modal inputs required. Generally, our proposed DAMC achieves the highest performance across all input combinations and tasks, demonstrating

Method	MCUB-3	MCUB-4
ImageBind-LLM	32.95	32.93
X-InstructBLIP	29.30	27.94
Proj-only	44.15	43.00
NaiveMC	54.70	54.03
DAMC	59.80	60.08

Table 3: Results on MCUB. MCUB-3 refers to subsets of the data with inputs from three modalities, while MCUB-4 includes inputs from four modalities.

its effectiveness of adeptly managing inputs from multiple modalities even without any training on these modal combinations. In addition, we make the following observations:

Model composition brings performance improvements with additional modal inputs. The results reveal a generally consistent trend where the application of model composition methods leads to improved performance as more modal inputs are integrated. For example, in Table 1 in the MUSIC-AVQA task, the transition from solely video (V) inputs to combinations that include image (I) and audio (A) inputs leads to noteworthy performance boosts. Specifically, when combining these three modalities, the performance for our DAMC reaches a peak at 57.32, marking a notable increase compared to V+I/A. In contrast, models like ImageBind-LLM and X-InstructBLIP do not show similar improvements when transitioning from V+I/A to the V+I+A combination. This demonstrates that models developed through model composition can handle inputs from different modalities more effectively, showcasing their ability to improve multimodal understanding.

DAMC outperforms strong baselines. Compared to previous methods, DAMC achieves superior performance, even surpassing methods that utilize paired multimodal data, such as ChatBridge and 3D-LLM (Hong et al., 2023). When comparing the performance across various model composition methods, DAMC significantly outperforms Proj-only and NaiveMC, particularly evident in scenarios involving the integration of more than two modalities. On the MCUB benchmark, for instance, DAMC achieves a significant performance enhancement, with an improvement over the second best NaiveMC of +5.10 points in scenarios combining three modalities and +6.05 points in the subset involving four modalities. This dis-



Figure 4: Qualitative examples on multimodal understanding of a composite model integrating four MLLMs.

inction proves the capability of DAMC to mitigate interference effectively when merging multiple MLLMs, demonstrating its robustness and efficiency for model composition.

4.4 Ablation Study

Our ablation study, summarized in Table 4, evaluates the effects of parameter decoupling and adjustment on DAMC. The findings reveal that employing neither strategy yields an average performance of 62.79 across benchmarks. Introducing parameter decoupling alone enhances the average to 65.05, while adjustment alone improves it to 63.55. Remarkably, combining both strategies boosts the average performance to 66.24, with significant improvements noted in all benchmarks. This highlights the synergistic impact of parameter decoupling and adjustment in optimizing multimodal model composition and achieving superior performance across varied tasks involving different types of modality inputs.

4.5 Qualitative Results

Figure 4 demonstrate the capability of the composition model to understand and reason over multimodality inputs. Additional qualitative results are in Appendix F.

Decoup.	Adjust.	MUSIC-AVQA	AVQA	MCUB-4	Avg.	Modal	NaiveMC	DAMC	Δ
\times	\times	53.63	80.70	54.03	62.79	I	1780.32	1772.85	-0.42%
\checkmark	\times	53.94	81.12	60.08	65.05	A + I	1718.42	1729.72	+0.66%
\times	\checkmark	55.36	81.27	54.03	63.55	V + I	1780.16	1825.50	+2.55%
\checkmark	\checkmark	57.32	81.31	60.08	66.24	V + I + A	1709.14	1852.48	+8.39%

Table 4: Ablation study on various components of DAMC. ‘‘Decoup.’’ and ‘‘Adjust.’’ denote parameter decoupling and adjustment.

Table 5: Composite model performance on the MME benchmark.

Parameter Decoupling		MUSIC-AVQA
V	I	V+I+A
No	No	55.36
Infer.	Train.	56.28
Train.	Infer.	56.42
Infer.	Infer.	56.16
Train.	Train.	57.32

Table 6: Results of different parameter decoupling strategies. ‘‘Infer.’’ and ‘‘Train.’’ denote parameter decoupling at inference or training time.

5 Discussion and Analysis

Composition of MLLMs with different training strategies.

In Section 3.3, we advocate training MLLMs with parameter decoupling. Intriguingly, this decoupling can also be implemented at inference time, by replicating the parameters to mimicking identical parameters on modality-specific and textual inputs. In this way, we can compose models with different training strategies. We conduct experiments on the MUSIC-AVQA dataset with V + I + A inputs. We vary the training strategy of video and image MLLMs, and compose them with an audio MLLMs trained with parameter decoupling. The results in Table 6 indicate that decoupling at inference time is also effective, which is beneficial when aiming to integrate existing models that were not trained with parameter decoupling.

Composite model performance on the single-modality task.

To explore the effectiveness of the composite model on single-modality tasks, we conducted experiments on the vision-language understanding evaluation benchmark MME (Fu et al., 2023). The results are shown in Table 5. Each composite model processes only images and textual questions. The Δ values denote the relative improvements of DAMC over NaiveMC. As more models are integrated, the performance of

NaiveMC tends to decline, while DAMC consistently shows improvement. This highlights the capability of our proposed methods to mitigate parameter interference effectively when composing models. Moreover, it is noteworthy that model composition enhances performance even in tasks limited to a single modality. We attribute this to the possibility that our approach could integrate knowledge from different models during the combination process, which we will leave for future research.

Impact of trainable parameters. We investigate whether the benefits of parameter decoupling arise from the augmentation of trainable parameters, given that two distinct sets of parameters are allocated for varying inputs. To this end, we train a new vision LLM without decoupling but increase the rank of the LoRA modules from 128 to 256, which doubles the trainable parameters. The findings, as shown in Table 7, reveal that merely expanding the pool of trainable parameters does not enhance performance. This suggests that the advantage of parameter decoupling extends beyond simple parameter quantity increase, highlighting its effectiveness in mitigating parameter interference.

Method	MUSIC-AVQA	AVQA
DAMC (r=128)	57.32	81.31
NaiveMC (r=128)	53.63	80.70
NaiveMC (r=256)	52.79	80.81

Table 7: Composite model performance on zero-shot audio-visual question answering tasks.

6 Conclusion

In this paper, we introduce and explore the concept of model composition for MLLMs, showcasing its ability to seamlessly integrate diverse modal capabilities without additional training. Initially, we introduce a foundational model composition framework, referred to as NaiveMC. Advancing

further, we design DAMC that employs parameter decoupling and adjustment to reduce cross-modal interference and optimize performance of the composite model. In addition, we construct a benchmark MCUB for multimodal commonality understanding evaluation to facilitate the relative research. Extensive experiments demonstrate the effectiveness of our approaches. We hope that our work will inspire further exploration into model composition for multimodal models.

Limitations

Our exploration is restricted to four commonly used modalities, omitting a comprehensive examination across the entire spectrum of potential modalities. This limitation may result in missed opportunities to identify further benefits or challenges associated with model composition. Additionally, our proposed approach has been tested on models of specific sizes, leaving the applicability of the methods on larger-scale MLLMs an open question for future research.

Acknowledgments

This work is supported by the National Key R&D Program of China (2022ZD0160502) and the National Natural Science Foundation of China (No. 61925601, 62276152).

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Feilong Chen, Minglun Han, Haozhi Zhao, Qingyang Zhang, Jing Shi, Shuang Xu, and Bo Xu. 2023. *X-llm: Bootstrapping advanced large language models by treating multi-modalities as foreign languages*.
- Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, and Furu Wei. 2022. Beats: Audio pre-training with acoustic tokenizers. *arXiv preprint arXiv:2212.09058*.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. *InstructBLIP: Towards general-purpose vision-language models with instruction tuning*. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. 2023. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153.
- Soham Deshmukh, Benjamin Elizalde, Rita Singh, and Huaming Wang. 2023. Pengi: An audio language model for audio tasks. *arXiv preprint arXiv:2305.11834*.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. 2023. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*.
- Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Manant Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. 2023. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15180–15190.
- Yuan Gong, Hongyin Luo, Alexander H Liu, Leonid Karlinsky, and James Glass. 2023. Listen, think, and understand. *arXiv preprint arXiv:2305.10790*.
- Jiaming Han, Renrui Zhang, Wenqi Shao, Peng Gao, Peng Xu, Han Xiao, Kaipeng Zhang, Chris Liu, Song Wen, Ziyu Guo, et al. 2023. Imagebind-llm: Multi-modality instruction tuning. *arXiv preprint arXiv:2309.03905*.
- Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 2023. 3d-llm: Injecting the 3d world into large language models. *arXiv preprint arXiv:2307.12981*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. 2023. Lorahub: Efficient cross-task generalization via dynamic lora composition. *arXiv preprint arXiv:2307.13269*.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2022. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*.
- Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. 2023. *Dataless knowledge fusion by merging weights of language models*.
- Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. 2019. Audiocaps: Generating captions for audios in the wild. In *NAACL-HLT*.
- Guangyao Li, Yake Wei, Yapeng Tian, Chenliang Xu, Ji-Rong Wen, and Di Hu. 2022. Learning to answer questions in dynamic audio-visual scenarios. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19108–19118.

- Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. 2023. Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*.
- Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. 2023. Scalable 3d captioning with pre-trained models.
- Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. 2023. Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration. *arXiv preprint arXiv:2306.09093*.
- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2023. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*.
- Michael Matena and Colin Raffel. 2022. Merging models with fisher-weighted averaging.
- Xinhao Mei, Chutong Meng, Haohe Liu, Qiuqiang Kong, Tom Ko, Chengqi Zhao, Mark D Plumbley, Yuexian Zou, and Wenwu Wang. 2023. Wavcaps: A chatgpt-assisted weakly-labelled audio captioning dataset for audio-language multimodal research. *arXiv preprint arXiv:2303.17395*.
- Artemis Panagopoulou, Le Xue, Ning Yu, Junnan Li, Dongxu Li, Shafiq Joty, Ran Xu, Silvio Savarese, Caiming Xiong, and Juan Carlos Niebles. 2023. X-instructblip: A framework for aligning x-modal instruction-aware representations to llms and emergent cross-modal reasoning. *arXiv preprint arXiv:2311.18799*.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks.
- Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. 2023. Pandagpt: One model to instruction-follow them all. *arXiv preprint arXiv:2305.16355*.
- Anirudh S Sundar, Chao-Han Huck Yang, David M Chan, Shalini Ghosh, Venkatesh Ravichandran, and Phani Sankar Nidadavolu. 2023. Multimodal attention merging for improved speech recognition and audio event classification. *arXiv preprint arXiv:2312.14378*.
- Yi-Lin Sung, Linjie Li, Kevin Lin, Zhe Gan, Mohit Bansal, and Lijuan Wang. 2023. An empirical study of multimodal model merging. *arXiv preprint arXiv:2304.14933*.
- Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, and Chao Zhang. 2024. SALMONN: Towards generic hearing abilities for large language models. In *The Twelfth International Conference on Learning Representations*.
- Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo-Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, and Ludwig Schmidt. 2022. Robust fine-tuning of zero-shot models.
- Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2023. Next-gpt: Any-to-any multimodal llm.
- Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 2015. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920.
- 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)*, pages 5288–5296.
- Runsen Xu, Xiaolong Wang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. 2023. Pointllm: Empowering large language models to understand point clouds. *arXiv preprint arXiv:2308.16911*.
- Pinci Yang, Xin Wang, Xuguang Duan, Hong Chen, Runze Hou, Cong Jin, and Wenwu Zhu. 2022. Avqa: A dataset for audio-visual question answering on videos. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 3480–3491.
- Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. 2024. Language models are super mario: Absorbing abilities from homologous models as a free lunch.
- Yan Zeng, Hanbo Zhang, Jiani Zheng, Jiangnan Xia, Guoqiang Wei, Yang Wei, Yuchen Zhang, and Tao Kong. 2023. What matters in training a gpt4-style language model with multimodal inputs? *arXiv preprint arXiv:2307.02469*.
- Hang Zhang, Xin Li, and Lidong Bing. 2023. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*.

Zijia Zhao, Longteng Guo, Tongtian Yue, Sihan Chen, Shuai Shao, Xinxin Zhu, Zehuan Yuan, and Jing Liu. 2023. Chatbridge: Bridging modalities with large language model as a language catalyst. *arXiv preprint arXiv:2305.16103*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#).

Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiayi Cui, HongFa Wang, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, et al. 2023. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment. *arXiv preprint arXiv:2310.01852*.

A Implementation Details of Pre-trained MLLMs

Table 9 details the components and training data of each MLLM across the four modalities, with specific explanations provided below:

- **Image:** We follow LLaVA-1.5 (Liu et al., 2023a) to use CLIP-ViT-L-336px as the image encoder with an MLP projection as the connector. The models are trained in a two-stage manner with LCS 558K (Liu et al., 2023b) as stage-1 data and LLaVA-mixed 665K (Liu et al., 2023a) as stage-2 data.
- **Audio:** We use BEATs-Iter3+ (Chen et al., 2022) as the audio encoder and a Q-Former with 32 query tokens as the connector following X-InstructBLIP (Panagopoulou et al., 2023). We use WaveCaps (Mei et al., 2023) for stage-1. For stage-2, we use a filtered version of OpenQA (Gong et al., 2023) with 350K examples.
- **Video:** Following Video-LLaVA (Lin et al., 2023), we use LanguageBind (Zhu et al., 2023) as the video encoder with an MLP connector. We reuse the stage-1 weights from Video-LLaVA for resource saving. The stage-2 data comprises Video-ChatGPT (Maaz et al., 2023) and a subset of LLaVA-mixed 665K including 100K image-text and 40K text-only instruction data.
- **Point cloud:** We use the pre-trained point encoder and instruction-following data from PointLLM (Xu et al., 2023). We use an MLP projection as the connector.

We adopt the same hyperparameters mainly following previous works (Liu et al., 2023a; Panagopoulou et al., 2023; Lin et al., 2023; Xu et al., 2023), as listed in Table 10. For the first training stage, only the parameters in the connectors are trainable. During the second training stage, for DAMC, since we decouple the parameters of modality inputs and text inputs, we additionally adjust the learning rate for the text components to $2e-5$ for all modalities. We apply the LoRA across all linear modules within the LLM, setting the LoRA rank to 128 and the alpha parameter to 256. For efficiency in training, we utilize DeepSpeed Zero Optimization stage 3.

Task	Image	Audio	Video	PC
MUSIC-AVQA	1	1/3	2/3	-
AVQA	1/3	2/3	2/3	-
MCUB-4	1/4	1/4	1/4	1/4

Table 8: Parameter adjustment weights for different MLLMs.

B Parameter Adjustment

In preliminary experiments, we find that composing MLLMs from two modalities typically achieved optimal results through a direct average, specifically a $1/2 + 1/2$ combination. Due to time and resource constraints, we only conducted parameter adjustments on three types of tasks. This process involved conducting a search for the optimal λ_i values within the set $[1/N, 2/N, \dots, N/N]$ for the composition of N modalities, based on the validation set performance for each corresponding task. The results are showcased in Table 8. We assume that the variance in coefficients comes from the differing demands of each task on the understanding capabilities across modalities. For MCUB-4, which requires a comprehensive grasp of content from all four modalities, an average coefficient emerged as the best result. Based on these findings, we also applied average coefficients for all MCUB-3 tasks.

C Details of Multimodal Commonality Understanding Benchmark

Table 11 presents the captioning datasets we use to generate MCUB task. For point cloud modality, we reserve 3000 point clouds from Objaverse (Deitke et al., 2023) dataset following Xu et al. (2023), and obtain their captions in Cap3D (Luo et al., 2023) dataset. Note that this part of the point clouds are **not** used in training and 3D object classification.

The semantic similarity of a group of entities is obtained by averaging the similarities of the captions of all two-entities combination, which are calculated by the all-MiniLM-L6-v2 model (Reimers and Gurevych, 2019). For example, if a group of entities with captions pair (A, B, C) is provided, the group similarity of it will be the average of the similarities between (A, B) , (A, C) and (B, C) .

We report the detailed results on subtasks of MCUB-3 benchmark in Table 12. The final result of MCUB-3 is reported as the average of the four sub-tasks.

Prompt template for generate questions, options

Modal	Modal Encoder	Connector	Stage-1 Data	Stage-2 Data
Image	CLIP-ViT-L-336px	MLP	LCS 558K	LLaVA-mixed 665K
Audio	BEATs-Iter3+	Q-Former	WaveCaps 400K	OpenQA filtered 350K
Video	LanguageBind	MLP	LCS 558K, Valley 702K	Video-ChatGPT 100K, LLaVA-mixed sampled 140K
Point Cloud	Point Encoder	MLP	PointLLM brief description 660K	Point complex instruction 70K

Table 9: Components and training data of MLLMs for different modalities.

Stage	Hyperparameter	Image	Audio	Video	Point Cloud
State-1	Batch size	256	256	256	128
	LR	1e-3	1e-3	1e-3	2e-3
	LR Schedule	cosine decay			
	Warmup Ratio	0.03			
	Epoch	1	1	1	3
Stage-2	Batch size	128	128	128	32
	LR	2e-4	1e-4	2e-4	2e-5
	LR Schedule	cosine decay			
	Warmup Ratio	0.03			
	Epoch	1	1	1	3

Table 10: Hyperparameters of different MLLMs.

Modal	Captioning Dataset
Image	COCO2017 (Lin et al., 2014) val set
Video	MSRVTT (Xu et al., 2016) test set
Audio	AudioCaps (Kim et al., 2019) test set
Point Cloud	Cap3D (Luo et al., 2023) (3000 subset)

Table 11: Captioning datasets in each modality to generate MCUB benchmark.

and correct answers:

Prompt Template

Given entity A with caption "A cat meowing and humans speaking on the background." with properties: agile, independent, and domesticated, entity B with caption "Loud barking and traffic" with properties aggressive, loud, and high energy, entity C with caption "Birds chirping in a quiet forest" and properties peaceful, wild, and vocal, and entity D with caption "A bustling city street with honking cars" and properties busy, noisy, and chaotic, you can generate a set of instruction answer pairs to find the common point of the entities as follows:

Example: Question: Which of the followings are the common point of the four entities.

A. Rhythmic ocean waves crashing on the shore. B. Gentle rustling of leaves in a serene garden. C. Audible environmental sounds. D. Soft crackling of a campfire under a starry sky. Answer: C. Explanation: The maximum common point among the four entities is the presence of ambient environmental sounds, which can be perceived audibly. Generate three such Question, Answer, Explanation triplets for entity A with caption "<Caption A>" and properties "<Properties A>", entity B with caption "<Caption B>" and properties "<Properties B>", entity C with caption "<Caption C>" and properties "<Properties C>", and entity D with caption "<Caption D>" and properties "<Properties D>" Examples:

D Prompt for Evaluation

We list the evaluation prompts for each dataset and modality combination in Table 13. In the prompts, we use "<image>", "<audio>", "<video>" and "<point>" to represent image, audio, video and point cloud modality inputs.

Method	MCUB-3			
	V+I+A	V+A+P	V+I+P	I+A+P
Imagebind-LLM	35.20	31.40	31.80	33.40
X-InstructBLIP	41.40	25.20	29.40	21.20
Proj-only	47.40	43.80	42.60	42.80
NaiveMC	56.00	51.00	53.00	58.80
DAMC	56.60	58.80	58.20	65.60

Table 12: Detailed results on sub-tasks of MCUB-3 involving different modalities combinations.

E Additional Point Cloud Results

In Table 2, we report exclusively on the zero-shot 3D object classification performance using P + I (Point cloud + Image) inputs, with comprehensive results detailed in Table 14. We find that ImageBind-LLM, X-InstructBLIP, and Proj-only struggle with point cloud inputs alone, particularly in adhering to open-ended generation instructions, leading to poor results. We assume that this issue likely stems from the necessity for point MLLM to undergo training on point-text instruction data with trainable LLM parameters to enhance performance.

F Additional Qualitative Results

We provide additional qualitative results about in Figure 5-7.

Dataset	Modal	Prompt Template
MUSIC-AVQA & AVQA	V	<video>\n{ Question} \nAnswer the question using a single word.
	V+I	Based on the video <video> and image <image>\n{ Question} \nAnswer the question using a single word.
	V+A	Based on the video <video> and audio <audio>\n{ Question} \nAnswer the question using a single word.
	V+I+A	Video: <video>\n Image: <image>\n Audio: <audio>\n { Question} \nAnswer the question using a single word.
Objaverse & ModelNet40	P	<point>\nWhat is this? (Type-I) / This is an object of (Type-C)
	P+I	Based on rendered image <image> and point cloud <point>\nWhat is this? (Type-I) / This is an object of (Type-C)
MCUB-3	V+I+A	Based on four input entities:\nimage <image>\naudio <audio>\nvideo <video>\n { Question} { Options} Answer with the option’s letter from the given choices directly.
	V+A+P	Based on four input entities:\naudio <audio>\nvideo <video>\npoint <point>\n { Question} { Options} Answer with the option’s letter from the given choices directly.
	V+I+P	Based on four input entities:\nimage <image>\nvideo <video>\npoint <point>\n { Question} { Options} Answer with the option’s letter from the given choices directly.
	I+A+P	Based on three input entities:\nimage <image>\naudio <audio>\npoint <point>\n { Question} { Options} Answer with the option’s letter from the given choices directly.
MCUB-4	V+I+A+P	Based on four input entities:\nimage <image>\naudio <audio>\nvideo <video>\npoint <point>\n { Question} { Options} Answer with the option’s letter from the given choices directly.

Table 13: Prompt Template for different evaluation benchmarks.

Task	Method	Instruction-type		Completion-type	
		P	P + I	P	P + I
Objaverse	3D-LLM	-	49.00	-	41.50
	Point-LLM	55.00	-	51.00	-
	ImageBind-LLM	*	31.00	*	26.50
	X-InstructBLIP	*	50.00	*	31.50
	Proj-only	17.50*	48.00	15.00*	42.50
	NaiveMC	56.00	55.00	54.50	59.50
	DAMC	57.00	60.50	56.50	62.00
ModelNet40	Point-LLM	53.44	-	51.82	-
	ImageBind-LLM	*	42.71	*	42.46
	X-InstructBLIP	*	61.43	*	61.14
	Proj-only	6.52*	62.88	5.92*	61.99
	NaiveMC	68.59	66.00	58.79	64.59
	DAMC	68.76	70.02	58.71	65.24

Table 14: Experimental results on zero-shot 3D object classification tasks using point (P) and combined point and image (P + I) inputs. Following Xu et al. (2023), two different types of prompts are considered: Instruction-type (“What is this?”) and Completion-type (“This is an object of”). *: models failed to follow open-ended generation instructions, leading to particularly low scores.

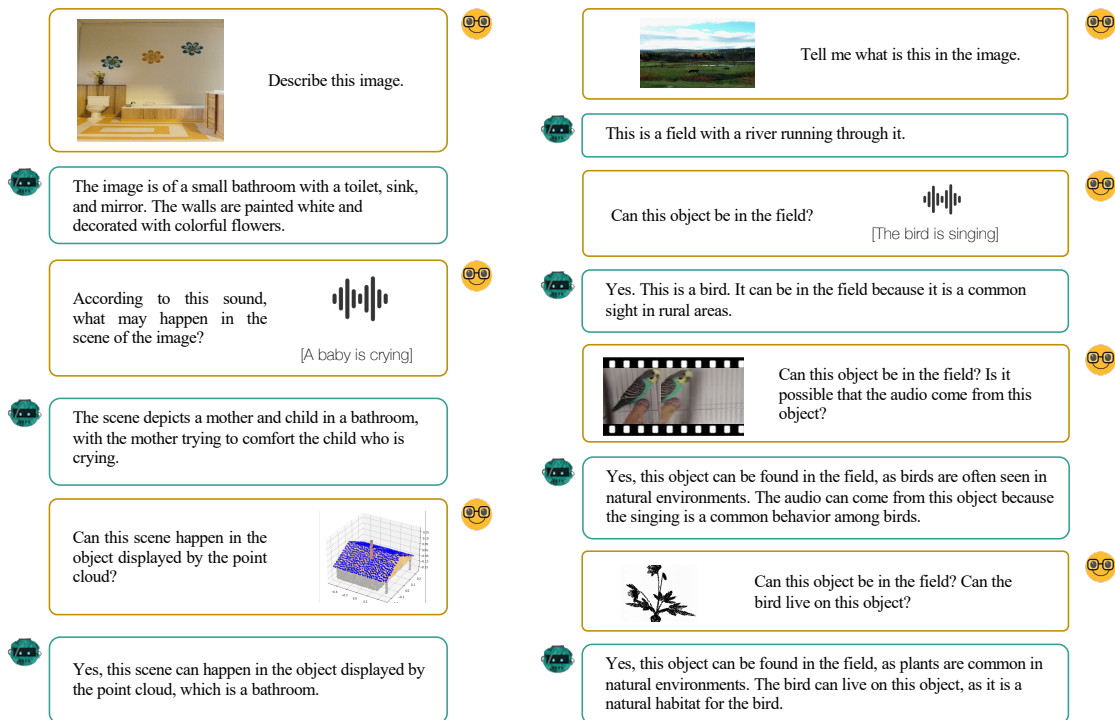




Figure 5: Additional qualitative results.



Describe this image in detail.

The image features a cat lying on a bed, with its head peeking out from under a blanket. The cat has a distinctive orange and white fur pattern, and its eyes are wide open, giving it a curious and alert appearance. The bed is covered with a blue blanket, which contrasts with the cat's orange and white fur. The cat's posture suggests that it is comfortable and relaxed, possibly enjoying the warmth and coziness of the bed.


Write a poem from this audio.



[birds squawk]

A flock of birds, so many and so bright,
Chirping and singing, their melodies delight.
Their wings flutter, their feathers they spread,
In the sky above, they dance and tread.

Their calls echo, their voices so clear,
A symphony of nature, so pure and dear.




How many whales are in this video?

There are two whales in this video.

What are they doing?

The two whales are swimming and splashing in the water.



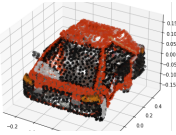

What is this species?

This is a dog.

Tell me more about this dog.

This dog is brown in color and has a long tail.



Figure 6: Additional qualitative results.

[The sound of car]

Can the object in the point cloud produce the sound in the audio?

Yes, the object in the point cloud is a car, and it is likely that the sound in the audio is produced by the car.

[The dog is snoring]

Can this sound come from the image? Why?

Yes, the sound of a dog snoring can come from the image. The dog is shown lying down and snoring, which is a common behavior of dogs.

Figure 7: Additional qualitative results.