

MACAW-LLM: MULTI-MODAL LANGUAGE MODELING WITH IMAGE, AUDIO, VIDEO, AND TEXT INTEGRATION

Chenyang Lyu^{1,2}, Minghao Wu³, Longyue Wang^{1*}, Xinting Huang¹,
Bingshuai Liu¹, Zefeng Du¹, Shuming Shi¹ & Zhaopeng Tu¹

¹Tencent AI Lab ²Dublin City University ³Monash University
chenyang.lyu2@mail.dcu.ie, minghao.wu@monash.edu,
{timxthuang,bsliu,zefengdu, shumingshi, zptu}@tencent.com

ABSTRACT

Although instruction-tuned large language models (LLMs) have exhibited remarkable capabilities across various NLP tasks, their effectiveness on other data modalities beyond text has not been fully studied. In this work, we propose MACAW-LLM, a novel multi-modal LLM that seamlessly integrates visual, audio, and textual information. MACAW-LLM consists of three main components: a modality module for encoding multi-modal data, a cognitive module for harnessing pretrained LLMs, and an *alignment module* for harmonizing diverse representations. Our novel alignment module seamlessly bridges multi-modal features to textual features, simplifying the adaptation process from the modality modules to the cognitive module. In addition, we construct a large-scale multi-modal instruction dataset in terms of multi-turn dialogue, including 69K image instances and 50K video instances. We have made our data, code and model publicly available, which we hope can pave the way for future research in multi-modal LLMs and expand the capabilities of LLMs to handle diverse data modalities and address complex real-world scenarios.



<https://github.com/lyuchenyang/Macaw-LLM>

1 INTRODUCTION

Instruction-tuned large language models (LLMs) have demonstrated impressive capabilities across various domains, exhibiting zero-shot generalization without the need for task-specific fine-tuning (Ouyang et al., 2022; Wei et al., 2022; Sanh et al., 2022; Chung et al., 2022; OpenAI, 2023). However, these models are primarily limited to processing text-based data. Previous research on multi-modal pre-training has shown promise in aligning knowledge from different modalities within a shared latent space (Wang et al., 2022a; Alayrac et al., 2022; Bao et al., 2022; Wang et al., 2022b). Furthermore, there is a recent line of research papers focusing on enabling multi-modal pre-trained models to understand and follow instructions (Xu et al., 2022; Zhu et al., 2023; Liu et al., 2023; Li et al., 2023a; Gong et al., 2023; Dai et al., 2023; Su et al., 2023; Huang et al., 2023).

In this work, we propose MACAW-LLM, a multi-modal instruction-tuned LLM that integrates four different modalities, including image, video, audio, and text, into one single model. We propose a novel alignment approach that aligns multi-modal features to the embeddings of LLMs, which produces aligned features that are closer to the textual features of language models and can be naturally injected into the input sequence of LLMs. A key motivation for our approach is to streamline the adaptation process for LLMs. In particular, MACAW-LLM employs a one-stage

*Longyue Wang is the corresponding author: vinnlywang@tencent.com.

arXiv:2306.09093v1 [cs.CL] 15 Jun 2023

instruction fine-tuning process, promoting a simpler learning experience. Previous multi-modal systems typically require two-stage training Li et al. (2023c); Zhu et al. (2023); Liu et al. (2023); Dai et al. (2023), where the first stage usually trains the projection layer for alignment between multi-modal features and text features, and the second stage is the general instruction fine-tuning for LLMs. In contrast, our approach aligns the multi-modal features to the embedding layer of LLMs, which produce aligned features based on LLMs embeddings that can be naturally injected into the input sequence of LLMs. This makes our approach more advantageous.

To address the limitations of current multi-modal datasets that predominantly emphasize specific task types, we create our MACAW-LLM instruction dataset, which is described in Section 4. This dataset covers a wide range of instructional tasks and combines various data modalities, making it more diverse and better-suited for multi-modal instruction-tuned LLMs. We utilize the remarkable generative capability of current LLMs, such as GPT-3.5-TURBO, to curate this dataset, ensuring the target text properly aligns with human instructions.

Our contributions in this work can be summarized as follows:

- We propose a novel architecture for multi-modal language modeling, which jointly learns to **align** multi-modal features and textual features and **generate** output sequence.
- We release MACAW-LLM instruction dataset, a large-scale **multi-modal instruction dataset** that covers diverse instructional tasks leveraging image and video modalities, which facilitates future work on multi-modal LLMs.

2 RELATED WORK

Instruction-Tuned Large Language Models Large language models (LLMs) have showcased exceptional generative capabilities in a wide range of natural language processing (NLP) tasks (Brown et al., 2020; Thoppilan et al., 2022; Hoffmann et al., 2022; Chowdhery et al., 2022). By leveraging techniques such as supervised instruction tuning and reinforcement learning from human feedback (RLHF), LLMs exhibit remarkable few- and zero-shot generalization capabilities (Ouyang et al., 2022; Wei et al., 2022; Sanh et al., 2022; Chung et al., 2022; Muennighoff et al., 2022; OpenAI, 2023; Anil et al., 2023). Recently, Wang et al. (2022c) highlight the lack of diversity in human-written instructions and demonstrate that machine-generated instructions can be used for instruction tuning. Since then, several instruction-tuned LLMs have been fine-tuned using various machine-generated instruction datasets (Taori et al., 2023; Chiang et al., 2023; Li et al., 2023b). More surprisingly, Wu et al. (2023b) reveal that instruction-following is not solely a property of LLMs, as even relatively small language models can follow instructions when fine-tuned on large-scale instruction datasets.

Multi-Modality Drawing inspiration from the human learning process, artificial intelligence (AI) researchers are actively exploring the combination of different modalities to train deep learning models. With the success of LLMs, feature alignment among multiple modalities has attracted great interest for its applications. There is a line of research works that learns a joint embedding space for multiple modalities (Radford et al., 2021; Baevski et al., 2022; Girdhar et al., 2023). Some researches also attempt to combine the pre-trained vision-only and language-only models, showcasing impressive zero-shot capabilities (Alayrac et al., 2022; Li et al., 2023c; Su et al., 2022). More recently, a number of works explore to enable the multi-modal LLMs to follow the instructions (Zhu et al., 2023; Ye et al., 2023; Li et al., 2023a; Chen et al., 2023; Gong et al., 2023; Dai et al., 2023). Xu et al. (2022) introduce MultiInstruct, the first multi-modal instruction tuning benchmark dataset covering a wide range of multi-modal tasks and categories. Liu et al. (2023) explore the multi-modal instruction-tuning using the machine-generated data. Su et al. (2023) allow the textual LLMs to support six modalities using the parameter-efficient fine-tuning technique LoRA.

Our Work In this work, we propose MACAW-LLM, a multi-modal LLM that effectively integrates information from visual, audio, and textual modalities, enabling it to comprehend and execute instructions accurately.

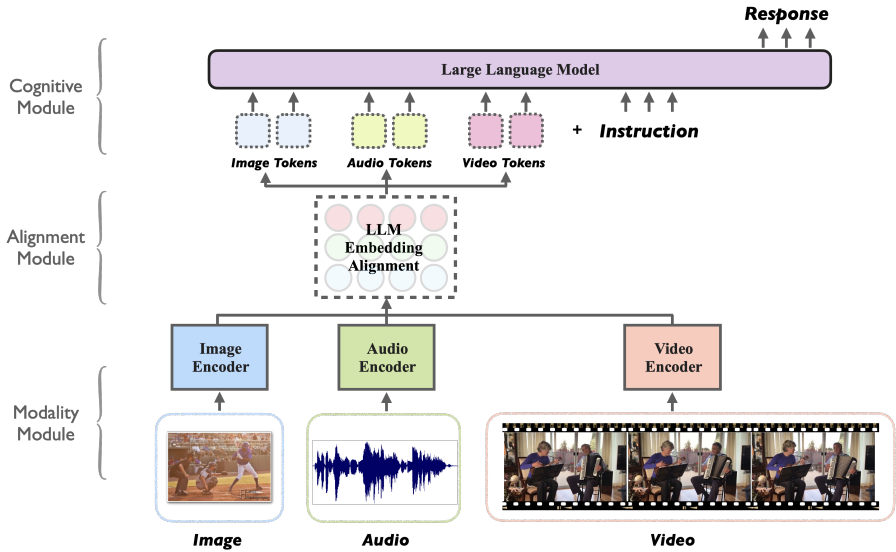


Figure 1: An overview of MACAW-LLM model architecture.

3 METHODOLOGY

In this section, we provide a comprehensive description of MACAW-LLM. We begin by presenting an outline of the model architecture, followed by a detailed description of each individual module within MACAW-LLM, namely the modality module, alignment module, and cognitive module. Lastly, we provide an in-depth explanation of the training process of MACAW-LLM.

3.1 MODEL OVERVIEW

We present an overview of MACAW-LLM in this section. As shown in Figure 1, there are three major modules in MACAW-LLM as follows:

- **Modality Module:** Existing LLMs primarily focus on processing textual information. To incorporate additional modalities such as visual and audio data, we integrate extra modality encoders into MACAW-LLM. This enhancement enables our MACAW-LLM to handle multiple modalities effectively.
- **Alignment Module:** Since each modality encoder is trained independently, the learned representations of different modalities may not be directly compatible. To address this, we propose the alignment module, which unifies the representations from different modalities, enabling effective integration of multi-modal information.
- **Cognitive Module:** LLMs have demonstrated remarkable capability in understanding and following human instructions. In MACAW-LLM, we leverage pretrained LLMs as our cognitive module, which forms the foundation of MACAW-LLM. It is worth noting that the cognitive module also serves as the textual modality encoder in our approach.

Figure 1 provides a visual representation of the MACAW-LLM architecture, while Section 3.2 and Section 3.3 offer detailed explanations of the modality module and alignment module, respectively. As the cognitive module of MACAW-LLM, the effectiveness of instruction-tuned LLMs has been demonstrated by several previous works (Ouyang et al., 2022; Wei et al., 2022; OpenAI, 2023; Taori et al., 2023; Chiang et al., 2023; Anil et al., 2023), and we follow their practices in our MACAW-LLM.

3.2 MODALITY MODULE

Existing LLMs are highly powerful but typically limited to processing only textual information. In this section, we describe how we encode information from different modalities.

Visual Modality Encoder Radford et al. (2021) propose a novel framework, known as CLIP (Radford et al., 2021), which exploits a significantly wider range of supervision by directly learning from unprocessed textual data related to images. In this work, we utilize the capabilities of CLIP-ViT-B/16 for encoding visual information, including images and video frames.

Audio Modality Encoder Radford et al. (2022) introduce a novel multilingual speech recognition model called WHISPER (Radford et al., 2022). This model is trained on a vast audio dataset with weak supervision. In MACAW-LLM, we leverage the power of WHISPER-BASE to encode the audio signals, thereby extracting meaningful representations from the audio data.

Textual Modality Encoder LLMs are commonly pre-trained on the massive text corpora, so instruction-tuned LLMs can naturally process text information. In this work, we consider LLAMA-7B (Touvron et al., 2023) as the foundation of MACAW-LLM.

We acknowledge the existence of numerous publicly available pre-trained models that can serve as modality encoders. However, we leave the investigation of their utility to the future work.

3.3 ALIGNMENT MODULE

Modality encoders are typically trained separately, leading to potential discrepancies in the representations generated by different encoders. As a result, it becomes crucial to align these independent representations within a joint space. In this section, we outline the approach we employ to align these representations.

Multi-Head Self-Attention (MHSA) Scaled dot-product attention is a fundamental component of the Transformer model (Vaswani et al., 2017). It operates on three inputs: the query vector $\mathbf{Q} \in \mathbb{R}^{n_q \times d_k}$, the key vector $\mathbf{K} \in \mathbb{R}^{n_k \times d_k}$, and the value vector $\mathbf{V} \in \mathbb{R}^{n_k \times d_v}$. This attention mechanism calculates attention weights by comparing the queries \mathbf{Q} with the keys \mathbf{K} . It then uses these weights to update the query representations through a weighted sum of the values \mathbf{V} and can be described as follows:

$$\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V}, \quad (1)$$

where d_k is the dimensionality of the key and query vectors, and n_q and n_k are the number of queries and keys, respectively.

Modality Alignment The alignment strategy is designed to efficiently connect multi-modal features with textual features, facilitating quicker adaptation. In this work, we refer to the image and video features obtained from our visual modality encoder (i.e. CLIP) as $\mathbf{x}_i \in \mathbb{R}^{L_i \times d_i}$ and $\mathbf{x}_v \in \mathbb{R}^{L_v \times d_v}$, respectively. Additionally, we denote the audio features from the audio modality encoder (i.e. WHISPER) as $\mathbf{x}_a \in \mathbb{R}^{L_a \times d_a}$. The process of modality alignment is outlined as follows:

1. **Encoding:** We firstly leverage the pre-trained models ,CLIP and WHISPER, to encode multi-modal features:

$$\mathbf{h}_i = \text{CLIP}(\mathbf{x}_i), \quad \mathbf{h}_v = \text{CLIP}(\mathbf{x}_v), \quad \mathbf{h}_a = \text{WHISPER}(\mathbf{x}_a), \quad (2)$$

where $\mathbf{h}_i \in \mathbb{R}^{L_i \times d_h}$, $\mathbf{h}_v \in \mathbb{R}^{L_v \times d_h}$ and $\mathbf{h}_a \in \mathbb{R}^{L_a \times d_h}$ are image, video, and audio features, respectively, and d_h is the dimension of modality-specific features.

2. **Transformation:** To reduce computational costs and minimize the number of tokens in the prefix, we employ a 1-D convolutional layer to compress the length of the multi-modal features to a smaller and fixed value. Subsequently, a linear layer is employed to adjust the hidden size of the features, aligning it with the size of the LLMs embeddings as follows:

$$\mathbf{h}'_i = \text{Linear}(\text{Conv1D}(\mathbf{h}_i)), \quad \mathbf{h}'_v = \text{Linear}(\text{Conv1D}(\mathbf{h}_v)), \quad \mathbf{h}'_a = \text{Linear}(\text{Conv1D}(\mathbf{h}_a)), \quad (3)$$

where $\mathbf{h}'_i \in \mathbb{R}^{L' \times d_e}$, $\mathbf{h}'_v \in \mathbb{R}^{L' \times d_e}$, and $\mathbf{h}'_a \in \mathbb{R}^{L' \times d_e}$ are the transformed features with a fixed length of L' and an embedding dimension of d_e . The value of L' is significantly smaller than L_i , L_v , and L_a , while d_e corresponds to the dimensionality of the embedding matrix $\mathbf{E} \in \mathbb{R}^{V \times d_e}$ associated with the textual LLMs (i.e. LLAMA in this work).

3. **Alignment:** Each modality encoder is trained separately, resulting in distinct representations for different modalities. To establish a common representation space, it becomes necessary to align these representations across modalities. In this work, we consider the transformed visual and audio modality representations obtained in Equation 3 as the *soft tokens* of LLM, the cognitive model, so we propose to align the visual and audio representations with the textual embedding space using the attention mechanism in Equation 1 as follows:

$$h^a = \text{Attn}(h', \mathbf{E}, \mathbf{E}), \tag{4}$$

where h' is the modality representation obtained in Equation 3 (i.e. h'_i, h'_v , and h'_a) and h^a is the corresponding aligned representation, specifically, h_i^a, h_v^a , and h_a^a . After such an alignment operation facilitated by the attention mechanism, the LLM (cognitive module) can seamlessly process the representations from various modalities.

4. **Integration:** The integration of aligned modality representations into the instruction can be achieved effortlessly through the concatenation operation. Given the aligned modality representations, the integration can be defined as follows:

$$\mathbf{x} = [h_i^a : h_v^a : h_a^a : \text{Embed}(x_t)], \tag{5}$$

where $[:]$ represents the concatenation operation, \mathbf{x} represents the multi-modal instruction, x_t represents the sequence of tokens in the textual instruction, and $\text{Embed}(x_t)$ represents the sequence of embeddings of x_t .

In this section, we describe how we align the multi-modality representation into a shared representation space using the attention mechanism. It is important to note that our model, MACAW-LLM, has the capability to process multiple modalities concurrently, while the textual instruction x_t is always necessary as part of the instruction \mathbf{x} . We intend to investigate the direct utilization of visual or audio instructions in our future work.

3.4 ONE-STEP INSTRUCTION FINE-TUNING

The common multi-modal practice in previous works involves two-step training (Li et al., 2023c; Liu et al., 2023; Dai et al., 2023). The first step focuses on training the projection layer to align multi-modal features with textual features, while the second step involves fine-tuning the general instruction for LLMs. In contrast, our approach, MACAW-LLM, simplifies the adaptation process by employing a one-step instruction fine-tuning approach. This approach ensures coherent alignment across the modalities and eliminates the potential risk of error propagation that can occur in multi-step fine-tuning procedures.

In this work, we fine-tune all the parameters θ in MACAW-LLM, and the objective is to minimize the negative log-likelihood over the response \mathbf{y} with respect to θ as follows:

$$\mathcal{L}(\mathbf{y}; \theta) = - \sum_{j=1}^N \log P(y_j | \mathbf{x}; \theta), \tag{6}$$

where N denotes the number of tokens in \mathbf{y} and y_j is the j -th token in \mathbf{y} . By employing such a one-step fine-tuning strategy, MACAW-LLM can effectively harmonize the different modules.

4 MACAW-LLM INSTRUCTION DATASET

Current multi-modal datasets, such as visual question answering (Antol et al., 2015; Goyal et al., 2017), summarization (Li et al., 2017; Jangra et al., 2023), and dialogue (Shuster et al., 2021; Sun et al., 2022), predominantly emphasize specific task types, resulting in a limited diversity of tasks. Additionally, the target text in these datasets often lacks proper alignment with the style of human-written text, making it difficult for models fine-tuned on such data to effectively follow human instructions. To address these limitations, we utilize the remarkable generative capability of current LLMs (i.e. GPT-3.5-TURBO) to curate our MACAW-LLM instruction dataset.

To generate the dataset, we utilize the power of GPT-3.5-TURBO. We provide it with a prompt in the form of an image or video caption (see Figure 3). To optimize the generation process and improve efficiency, we generate 10 instruction-response pairs within a single query. For image

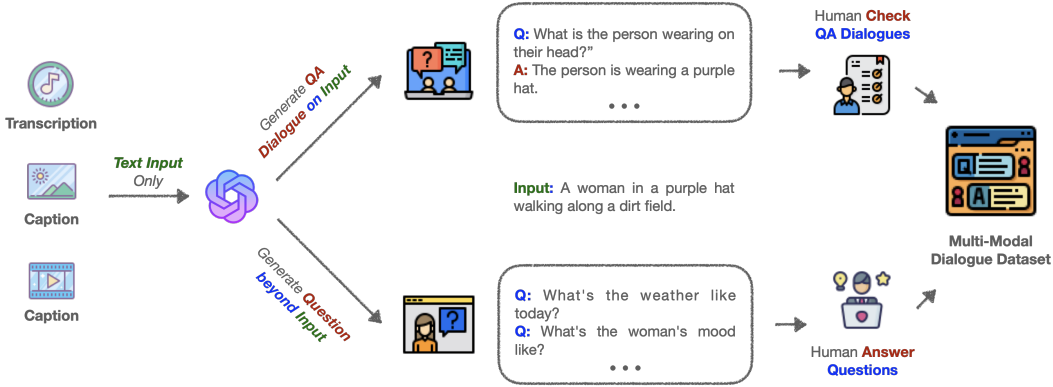


Figure 2: The process of constructing a multi-modal instruction corpus. GPT-4 is prompted to generate instruction-response examples from text input, which then correspond to its multi-modal contents. Human verification and annotation are incorporated to ensure the accuracy.

caption data, we rely on the MS COCO dataset (Lin et al., 2014). It consists of 328,000 images accompanied by captions. From this dataset, we randomly select a subset of 10,000 images with their respective captions to create our dataset. In addition to image data, we incorporate video caption data from two datasets: Charades (Sigurdsson et al., 2016) and AVSD (AlAmri et al., 2019). These datasets collectively contain 9,848 videos with captions, which we utilize to create our own dataset.

We repeat this process and obtain approximately 69K examples based on COCO image captions and about 50K examples based on Charades and AVSD video captions. The dataset creation process is illustrated in Figure 2. Table 1 provides statistics about the dataset, including the number of items, the word count of instructions and responses, and examples of each type.

Our current dataset is focused on single-turn dialogues, but we acknowledge the significance of including multi-turn dialogues and expanding the dataset to encompass a wider range of multi-modal content. To address this, we are actively engaged in the process of incorporating multi-turn dialogues and diversifying the dataset to enhance its richness. These additions will greatly contribute to enriching the dataset and will be beneficial for the fine-tuning process of language learning models (LLMs).

5 EXPERIMENTAL SETUP

5.1 DATASET

In this study, we utilize instruction data from three different sources:

- **Text instruction dataset:** For textual instruction-tuning, we make use of the Alpaca instruction dataset (Taori et al., 2023), which comprises approximately 52,000 instruction-response examples distilled from the TEXT-DAVINCI-003 model.
- **Image instruction dataset:** To create an image instruction dataset, we curate around 69K instruction-response pairs by generating them from COCO image captions (Lin et al., 2014) using GPT-3.5-TURBO as described in Section 4.
- **Video instruction data:** We generate approximately 50K video instruction-response examples by utilizing the video captions from the Charades (Sigurdsson et al., 2016) and AVSD (AlAmri et al., 2019) datasets using GPT-3.5-TURBO as described in Section 4.

In practice, we randomly sample 50K examples from each type of instruction data and combine them to form a final training dataset consisting of 150K examples. Note that the audio inputs are currently associated with the video instruction data and we are actively in the process of creating the audio instruction dataset.

Building Multi-Modal Dialogue with GPT-4

Prompt:

This is the **caption of an image/video: {...}**. This image/video contains important information that needs to be conveyed through high-quality instructions.

Your task is to provide ten pairs of instructions and responses that are related to the content of the image/video caption like dialogue concentrating on the content of the image/video without explicitly mentioning the caption or the word 'caption'.

Your focus should be on describing, explaining, or analyzing various aspects of the image/video, as well as providing some QA pairs. The purpose of this exercise is to fine-tune a language model so that it can generate accurate and relevant responses.

In each pair, the first line should start with "Q:" and contain an instruction related to the image/video, while the second line should start with "A:" and provide a response to the instruction.

Please ensure that your instructions are diverse and of high quality, accurately reflecting the content of the image and providing useful information to the language model:

Output:

Q: *Can you describe the color of the river in the image?*
A: **The river in the image appears to be a tranquil shade of blue.**

Q: *What type of boat is the man in the image paddling?*
A: **The man in the image is paddling a kayak.**

Q: *How do you think the man in the image is feeling while paddling down the river?*
A: **Judging by the peaceful surroundings and the calm pace of the paddling, it's likely that the man in the image is feeling relaxed and at ease.**

.....

Figure 3: Our instruction generation prompt for creating instruction-response pairs related to the content of an image/video caption using GPT-4. The objective is to create high-quality instructions and responses without explicitly mentioning the caption, aiming to improve the language model's ability to generate accurate and relevant responses.

5.2 HYPERPARAMETERS

We utilize DeepSpeed (Rasley et al., 2020) for optimization during the training process. The training is conducted on 8 Nvidia A100 GPUs. For each device, the training batch size is set to 4. We employ a gradient accumulation step of 3. The model is trained for 5 epochs, with a learning rate of 3×10^{-5} . The warmup ratio is 0.03, along with a cosine learning rate scheduler. The maximum sequence length is fixed at 512. We use FP16 precision for both training and inference.

6 EXAMPLES

To showcase the effectiveness and potential of our proposed MACAW-LLM in creating human-like conversational agents, this section provides compelling examples that demonstrate the system's remarkable ability to understand and generate responses related to visual content. These examples vividly illustrate how MACAW-LLM seamlessly processes and integrates multiple modalities of information, such as visuals and audio, within the domain of natural language processing (NLP). By generating informative, relevant, and coherent responses to a wide range of questions, MACAW-LLM clearly demonstrates its proficiency in NLP and underscores its potential for developing highly effective human-machine communication interfaces.



Dataset	Items	Ins. Len.	Res. Len.	Example
COCO	69,314	10.1	15.7	 <p>Caption: A woman is trying to hit a ball with a tennis racket. Instruction: What sport is the woman playing in the image? Response: The woman in the image is playing tennis with a racket.</p>
Charades/ AVSD	50,656	10.2	14.9	 <p>Caption: A person is standing next to a table holding a sandwich in their hand. They sit the sandwich on the table and pick up a book. Instruction: What does the person do right after setting down the sandwich? Response: The person picks up a book that is on the table.</p>

Table 1: Statistics and examples of our instruction dataset. We report the number of items, average instruction length, and average response length of our instruction dataset, based on COCO, Charades, and AVSD.

We present several examples that highlight the proficiency of our MACAW-LLM in understanding and following multi-modal instructions. In Figure 4, Figure 5, and Figure 6, we showcase our system’s multi-modal ability to understand and generate responses based on an image. These examples demonstrate how our system comprehends visual content and produces high-quality, fluent responses in natural language conversations. Our system generates contextually relevant and informative answers to various questions about the image, demonstrating its capability to communicate about visual content naturally and fluently. Figure 7 and Figure 8 present two examples that demonstrate MACAW-LLM’s excellent understanding of videos. We showcase its responses to various questions related to the video content, highlighting its ability to comprehend video information effectively. Furthermore, Figure 9 demonstrates our system’s capacity to process and integrate multiple modalities of information simultaneously. In this example, in addition to answering various video-grounded questions, MACAW-LLM effectively identifies whether the dog in the video is barking or not.

In summary, the examples provided showcase the impressive capabilities of our system in generating top-notch, contextually appropriate, and logically consistent responses to diverse questions about visual content within a natural language conversation. The proficiency of our system in natural language processing (NLP) and its adeptness in seamlessly incorporating multiple modalities of information underscore its tremendous potential in designing efficient interfaces for human-machine communication.

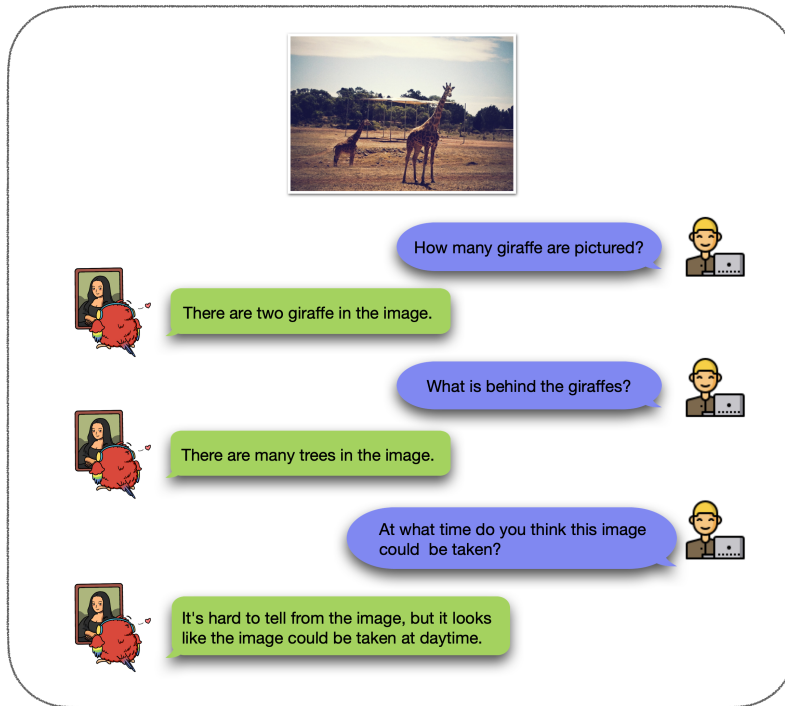


Figure 4: An example showcasing MACAW-LLM’s basic capability in **image-grounded question answering**. The image features two giraffes with a backdrop of numerous trees. MACAW-LLM can identifies these contents and infers that the photo was taken at daytime.

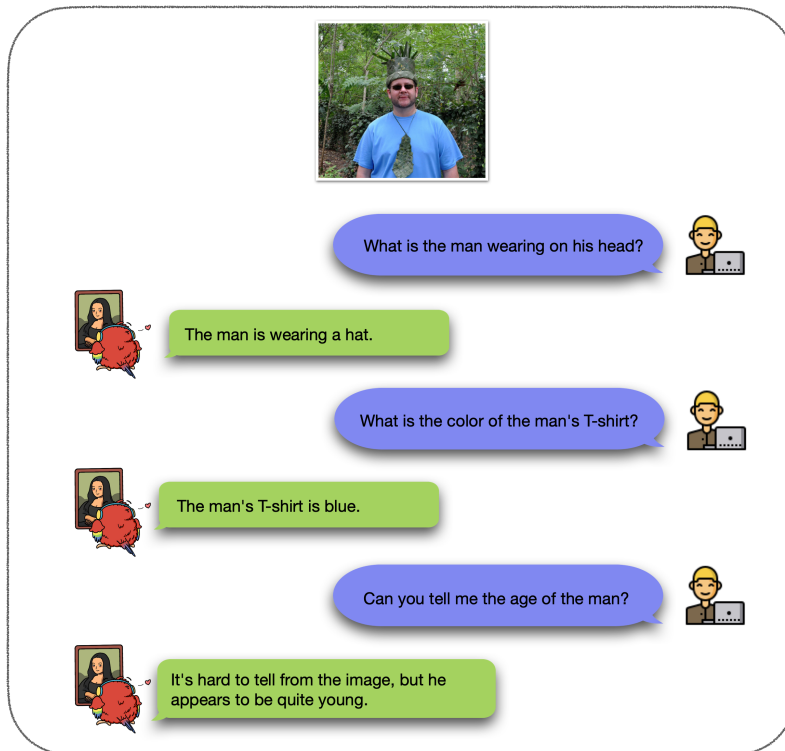


Figure 5: An example showcasing MACAW-LLM’s capability in **image-grounded understanding and reasoning**. As seen, MACAW-LLM can comprehend fundamental objects, such as a hat and a T-shirt. Besides, it tries to estimate the age of the man.

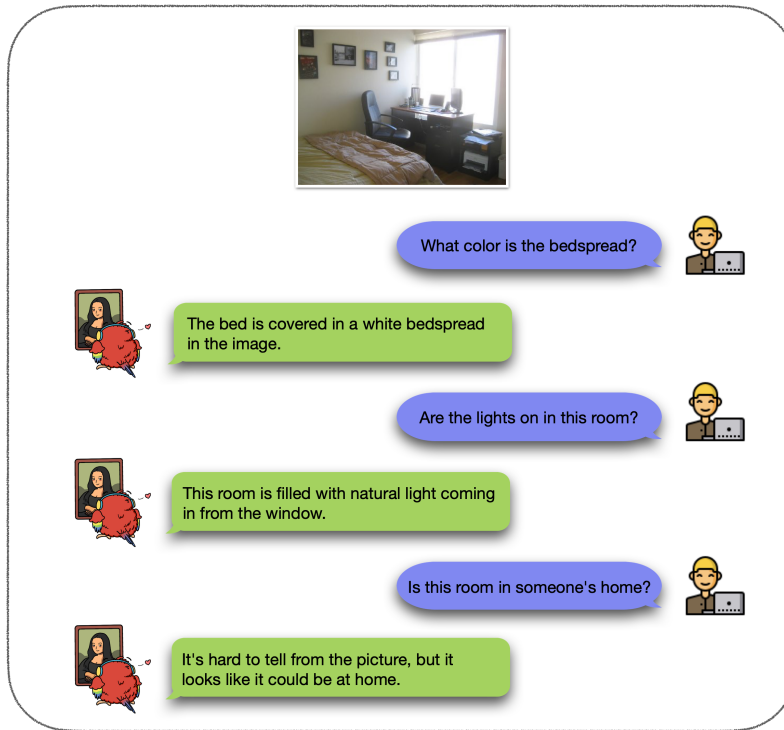


Figure 6: An example showcasing MACAW-LLM’s capability on **recognizing color and light**. Besides, MACAW-LLM estimate the location of the room.

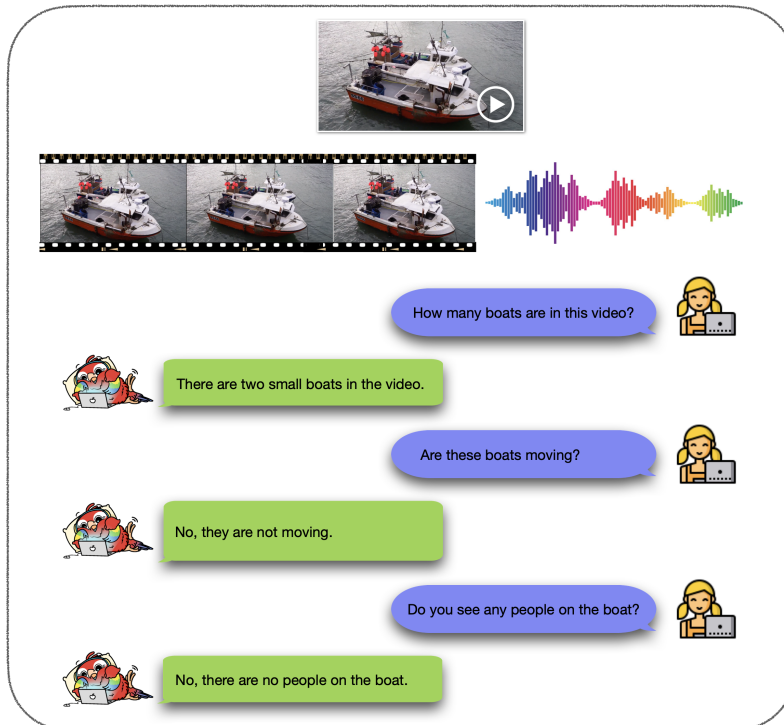


Figure 7: An example showcasing MACAW-LLM’s capability in **video-grounded question answering**. MACAW-LLM can recognize the boats and their amount. Besides, it is able to identify boats’ actions over time.

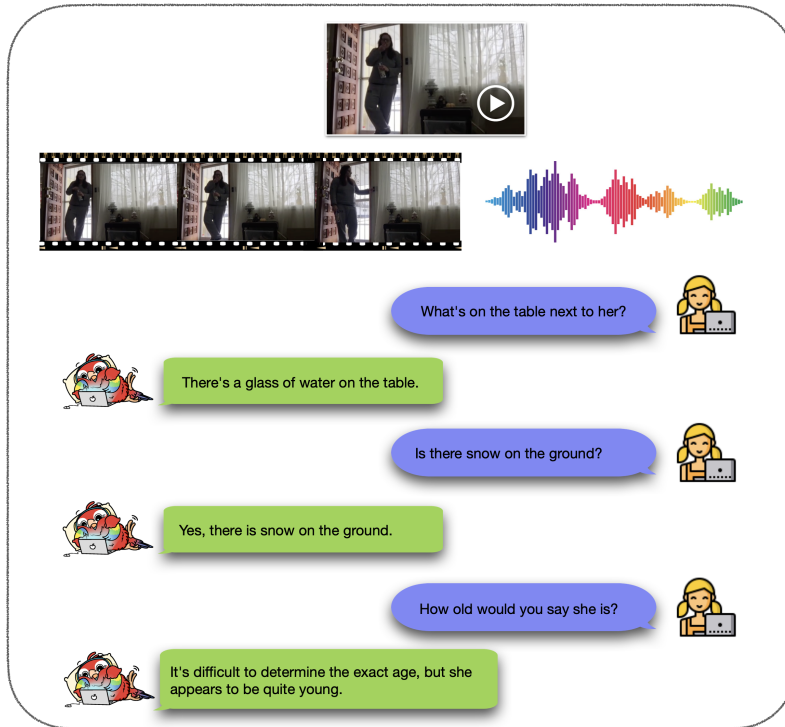


Figure 8: An example showcasing MACAW-LLM’s capability in **visual reasoning**. Despite only a small portion of “white” being visible outside the door, MACAW-LLM can infer the presence of “snow”. Furthermore, it attempts to estimate the age of the woman.



Figure 9: An example showcasing MACAW-LLM’s capability in **video- and audio-grounded question answering**. The video showcases a dog on a grassy field, remaining silent as indicated by the audio track.

7 LIMITATIONS

In this section, we summarize the limitations of MACAW-LLM as follows:

- **Evaluation:** We show some examples showcasing the multi-modal ability of our MACAW-LLM. However, we acknowledge that these efforts may not be fully adequate for accurately and comprehensively demonstrate model capabilities. Gudibande et al. (2023) highlights that instruction-tuned LLMs might not perform as well as the reported evaluation results suggest. Hence, we have concerns regarding the ability of our evaluation to provide an accurate reflection of the true capabilities of MACAW-LLM.
- **Single-Turn Dialogue:** While our training data mainly consists of "dialog-like" instructions, it's important to note that these instructions are currently limited to single-turn interactions. It is crucial to acknowledge that MACAW-LLM are not currently optimized for handling multi-turn dialogues and may not effectively leverage long-range context.
- **Hallucination, Toxicity and Fairness:** According to empirical evidence presented by Wu et al. (2023b), instruction-tuned LLMs may encounter issues such as hallucination, toxicity, and fairness. However, it is important to note that we do not evaluate our models, MACAW-LLM, in relation to these aspects due to the unavailability of suitable evaluation suites.

We acknowledge these limitations and recognize the need for addressing them in future work.

8 CONCLUSION AND FUTURE WORK

In this paper, we present MACAW-LLM, a multi-modal instruction-tuned LLM that accommodates four distinct modalities: image, video, audio, and text. In addition to the standard modality module and cognitive module, we propose a novel approach to align representations from different modality encoders into a shared space. Unlike previous methods, our approach combines representation alignment and instruction tuning into a single step, mitigating potential error propagation during multi-step tuning. Furthermore, we curate MACAW-LLM instruction dataset, a large-scale dataset of multi-modal instructions using GPT-3.5-TURBO. We demonstrate examples showcasing the multi-modal understanding ability of MACAW-LLM.

We discuss the limitations of our work and point out that current multi-modal instruction-tuned LLMs may suffer from various aspects in Section 7. We leave the investigation of these issues to the future work. Furthermore, we intend to broaden our corpus to encompass multi-turn and multilingual dialogues. This endeavor will take advantage of the capabilities of LLMs to effectively generate/translate long-document texts (Wang et al., 2017; Lyu et al., 2023; Wang et al., 2023; Wu et al., 2023a).

REFERENCES

- Huda AlAmri, Vincent Cartillier, Abhishek Das, Jue Wang, Anoop Cherian, Irfan Essa, Dhruv Batra, Tim K. Marks, Chiori Hori, Peter Anderson, Stefan Lee, and Devi Parikh. Audio visual scene-aware dialog. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 7558–7567. Computer Vision Foundation / IEEE, 2019. doi: 10.1109/CVPR.2019.00774. URL http://openaccess.thecvf.com/content_CVPR_2019/html/Alamri_Audio_Visual_Scene-Aware_Dialog_CVPR_2019_paper.html.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. Flamingo: a visual language model for few-shot learning. In *NeurIPS, 2022*. URL http://papers.nips.cc/paper_files/paper/2022/hash/960a172bc7fbf0177cccbb411a7d800-Abstract-Conference.html.
- Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H.

- Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernández Ábrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan A. Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vladimir Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, and et al. Palm 2 technical report. *CoRR*, abs/2305.10403, 2023. doi: 10.48550/arXiv.2305.10403. URL <https://doi.org/10.48550/arXiv.2305.10403>.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015. URL https://openaccess.thecvf.com/content_iccv_2015/html/Antol_VQA_Visual_Question_ICCV_2015_paper.html.
- Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. data2vec: A general framework for self-supervised learning in speech, vision and language. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 1298–1312. PMLR, 2022. URL <https://proceedings.mlr.press/v162/baevski22a.html>.
- Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. In *ICLR 2022*, April 2022. URL <https://www.microsoft.com/en-us/research/publication/beit-bert-pre-training-of-image-transformers/>.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf.
- Feilong Chen, Minglun Han, Haozhi Zhao, Qingyang Zhang, Jing Shi, Shuang Xu, and Bo Xu. X-LLM: bootstrapping advanced large language models by treating multi-modalities as foreign languages. *CoRR*, abs/2305.04160, 2023. doi: 10.48550/arXiv.2305.04160. URL <https://doi.org/10.48550/arXiv.2305.04160>.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://vicuna.lmsys.org>.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways. *CoRR*, abs/2204.02311, 2022. doi: 10.48550/arXiv.2204.02311. URL <https://doi.org/10.48550/arXiv.2204.02311>.

- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models. *CoRR*, abs/2210.11416, 2022. doi: 10.48550/arXiv.2210.11416. URL <https://doi.org/10.48550/arXiv.2210.11416>.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven C. H. Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *CoRR*, abs/2305.06500, 2023. doi: 10.48550/arXiv.2305.06500. URL <https://doi.org/10.48550/arXiv.2305.06500>.
- Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. *CoRR*, abs/2305.05665, 2023. doi: 10.48550/arXiv.2305.05665. URL <https://doi.org/10.48550/arXiv.2305.05665>.
- Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. Multimodal-gpt: A vision and language model for dialogue with humans. *CoRR*, abs/2305.04790, 2023. doi: 10.48550/arXiv.2305.04790. URL <https://doi.org/10.48550/arXiv.2305.04790>.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. URL https://openaccess.thecvf.com/content_cvpr_2017/html/Goyal_Making_the_v_CVPR_2017_paper.html.
- Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and Dawn Song. The false promise of imitating proprietary llms. *CoRR*, abs/2305.15717, 2023. doi: 10.48550/arXiv.2305.15717. URL <https://doi.org/10.48550/arXiv.2305.15717>.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models. *CoRR*, abs/2203.15556, 2022. doi: 10.48550/arXiv.2203.15556. URL <https://doi.org/10.48550/arXiv.2203.15556>.
- Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Barun Patra, Qiang Liu, Kriti Aggarwal, Zewen Chi, Johan Bjorck, Vishrav Chaudhary, Subhojit Som, Xia Song, and Furu Wei. Language is not all you need: Aligning perception with language models, 2023. URL <https://doi.org/10.48550/arXiv.2302.14045>.
- Anubhav Jangra, Sourajit Mukherjee, Adam Jatowt, Sriparna Saha, and Mohammad Hasanuzzaman. A survey on multi-modal summarization. *ACM Comput. Surv.*, feb 2023. ISSN 0360-0300. doi: 10.1145/3584700. URL <https://doi.org/10.1145/3584700>. Just Accepted.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *CoRR*, abs/2305.03726, 2023a. doi: 10.48550/arXiv.2305.03726. URL <https://doi.org/10.48550/arXiv.2305.03726>.
- Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji, and Timothy Baldwin. Bactrian-x : A multilingual replicable instruction-following model with low-rank adaptation. *CoRR*, abs/2305.15011, 2023b. doi: 10.48550/arXiv.2305.15011. URL <https://doi.org/10.48550/arXiv.2305.15011>.
- Haoran Li, Junnan Zhu, Cong Ma, Jiajun Zhang, and Chengqing Zong. Multi-modal summarization for asynchronous collection of text, image, audio and video. In *Proceedings of*

- the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 1092–1102, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1114. URL <https://aclanthology.org/D17-1114>.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. *CoRR*, abs/2301.12597, 2023c. doi: 10.48550/arXiv.2301.12597. URL <https://doi.org/10.48550/arXiv.2301.12597>.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In David J. Fleet, Tomás Pajdla, Bernt Schiele, and Tinne Tuytelaars (eds.), *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V*, volume 8693 of *Lecture Notes in Computer Science*, pp. 740–755. Springer, 2014. doi: 10.1007/978-3-319-10602-1_48. URL https://doi.org/10.1007/978-3-319-10602-1_48.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *CoRR*, abs/2304.08485, 2023. doi: 10.48550/arXiv.2304.08485. URL <https://doi.org/10.48550/arXiv.2304.08485>.
- Chenyang Lyu, Jitao Xu, and Longyue Wang. New trends in machine translation using large language models: Case examples with chatgpt. *CoRR*, abs/2305.01181, 2023. doi: 10.48550/arXiv.2305.01181. URL <https://doi.org/10.48550/arXiv.2305.01181>.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M. Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. Crosslingual generalization through multitask finetuning. *CoRR*, abs/2211.01786, 2022. doi: 10.48550/arXiv.2211.01786. URL <https://doi.org/10.48550/arXiv.2211.01786>.
- OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/arXiv.2303.08774. URL <https://doi.org/10.48550/arXiv.2303.08774>.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=TG8KACxEON>.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8748–8763. PMLR, 2021. URL <http://proceedings.mlr.press/v139/radford21a.html>.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. *CoRR*, abs/2212.04356, 2022. doi: 10.48550/arXiv.2212.04356. URL <https://doi.org/10.48550/arXiv.2212.04356>.
- Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In Rajesh Gupta, Yan Liu, Jiliang Tang, and B. Aditya Prakash (eds.), *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*, pp. 3505–3506. ACM, 2020. doi: 10.1145/3394486.3406703. URL <https://doi.org/10.1145/3394486.3406703>.

- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. Multitask prompted training enables zero-shot task generalization. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=9Vrb9D0WI4>.
- Kurt Shuster, Eric Michael Smith, Da Ju, and Jason Weston. Multi-modal open-domain dialogue. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 4863–4883, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.398. URL <https://aclanthology.org/2021.emnlp-main.398>.
- Gunnar A. Sigurdsson, Olga Russakovsky, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. Much ado about time: Exhaustive annotation of temporal data. *CoRR*, abs/1607.07429, 2016. URL <http://arxiv.org/abs/1607.07429>.
- Yixuan Su, Tian Lan, Yahui Liu, Fangyu Liu, Dani Yogatama, Yan Wang, Lingpeng Kong, and Nigel Collier. Language models can see: Plugging visual controls in text generation. *CoRR*, abs/2205.02655, 2022. doi: 10.48550/arXiv.2205.02655. URL <https://doi.org/10.48550/arXiv.2205.02655>.
- Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. Pandagpt: One model to instruction-follow them all. *CoRR*, abs/2305.16355, 2023. doi: 10.48550/arXiv.2305.16355. URL <https://doi.org/10.48550/arXiv.2305.16355>.
- Qingfeng Sun, Yujing Wang, Can Xu, Kai Zheng, Yaming Yang, Huang Hu, Fei Xu, Jessica Zhang, Xiubo Geng, and Daxin Jiang. Multimodal dialogue response generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2854–2866, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.204. URL <https://aclanthology.org/2022.acl-long.204>.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Kathleen S. Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Agüera-Arcas, Claire Cui, Marian Croak, Ed H. Chi, and Quoc Le. Llama: Language models for dialog applications. *CoRR*, abs/2201.08239, 2022. URL <https://arxiv.org/abs/2201.08239>.
- 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. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971, 2023. doi: 10.48550/arXiv.2302.13971. URL <https://doi.org/10.48550/arXiv.2302.13971>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von

- Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
- Longyue Wang, Zhaopeng Tu, Andy Way, and Qun Liu. Exploiting cross-sentence context for neural machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 2826–2831, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1301. URL <https://aclanthology.org/D17-1301>.
- Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. Document-level machine translation with large language models. *CoRR*, abs/2304.02210, 2023. doi: 10.48550/arXiv.2304.02210. URL <https://doi.org/10.48550/arXiv.2304.02210>.
- Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. OFA: unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 23318–23340. PMLR, 2022a. URL <https://proceedings.mlr.press/v162/wang22a1.html>.
- Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, and Furu Wei. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. *CoRR*, abs/2208.10442, 2022b. doi: 10.48550/arXiv.2208.10442. URL <https://doi.org/10.48550/arXiv.2208.10442>.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions. *CoRR*, abs/2212.10560, 2022c. doi: 10.48550/arXiv.2212.10560. URL <https://doi.org/10.48550/arXiv.2212.10560>.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations, 2022*. URL <https://openreview.net/forum?id=gEZrGCozdqR>.
- Minghao Wu, George Foster, Lizhen Qu, and Gholamreza Haffari. Document flattening: Beyond concatenating context for document-level neural machine translation. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pp. 448–462, Dubrovnik, Croatia, May 2023a. Association for Computational Linguistics. URL <https://aclanthology.org/2023.eacl-main.33>.
- Minghao Wu, Abdul Waheed, Chiyu Zhang, Muhammad Abdul-Mageed, and Alham Fikri Aji. Lamini-lm: A diverse herd of distilled models from large-scale instructions. *CoRR*, abs/2304.14402, 2023b. doi: 10.48550/arXiv.2304.14402. URL <https://doi.org/10.48550/arXiv.2304.14402>.
- Zhiyang Xu, Ying Shen, and Lifu Huang. Multiinstruct: Improving multi-modal zero-shot learning via instruction tuning. *CoRR*, abs/2212.10773, 2022. doi: 10.48550/arXiv.2212.10773. URL <https://doi.org/10.48550/arXiv.2212.10773>.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. mplug-owl: Modularization empowers large language models with multimodality. *CoRR*, abs/2304.14178, 2023. doi: 10.48550/arXiv.2304.14178. URL <https://doi.org/10.48550/arXiv.2304.14178>.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *CoRR*, abs/2304.10592, 2023. doi: 10.48550/arXiv.2304.10592. URL <https://doi.org/10.48550/arXiv.2304.10592>.