

CAT : Enhancing Multimodal Large Language Model to Answer Questions in Dynamic Audio-Visual Scenarios

Qilang Ye¹, Zitong Yu^{1*}, Rui Shao², Xinyu Xie¹,
Philip Torr³, and Xiaochun Cao⁴

¹ Great Bay University

² Harbin Institute of Technology, Shenzhen

³ University of Oxford

⁴ Shenzhen Campus of Sun Yat-sen University

Abstract. This paper focuses on the challenge of answering questions in scenarios that are composed of rich and complex dynamic audio-visual components. Although existing Multimodal Large Language Models (MLLMs) can respond to audio-visual content, these responses are sometimes ambiguous and fail to describe specific audio-visual events. To overcome this limitation, we introduce the CAT, which enhances MLLM in three ways: 1) besides straightforwardly bridging audio and video, we design a clue aggregator that aggregates question-related clues in dynamic audio-visual scenarios to enrich the detailed knowledge required for large language models. 2) CAT is trained on a mixed multimodal dataset, allowing direct application in audio-visual scenarios. Notably, we collect an audio-visual joint instruction dataset named AVInstruct, to further enhance the capacity of CAT to model cross-semantic correlations. 3) we propose AI-assisted ambiguity-aware direct preference optimization, a strategy specialized in retraining the model to favor the non-ambiguity response and improve the ability to localize specific audio-visual objects. Extensive experimental results demonstrate that CAT outperforms existing methods on multimodal tasks, especially in Audio-Visual Question Answering (AVQA) tasks. The codes and the collected instructions are released at <https://github.com/rikeilong/Bay-CAT>.

Keywords: Multimodal Large Language Model · Audio-visual Question Answering

1 Introduction

The real world revolves around sound and visual information, and their combination enhances our ability to perceive the world. Similarly, the development of Multimodal Large Language Models (MLLMs) [14, 64, 67] are closely related to audio and visual. Supervised fine-tuning [35, 46, 66] with specific instruction

* Corresponding author

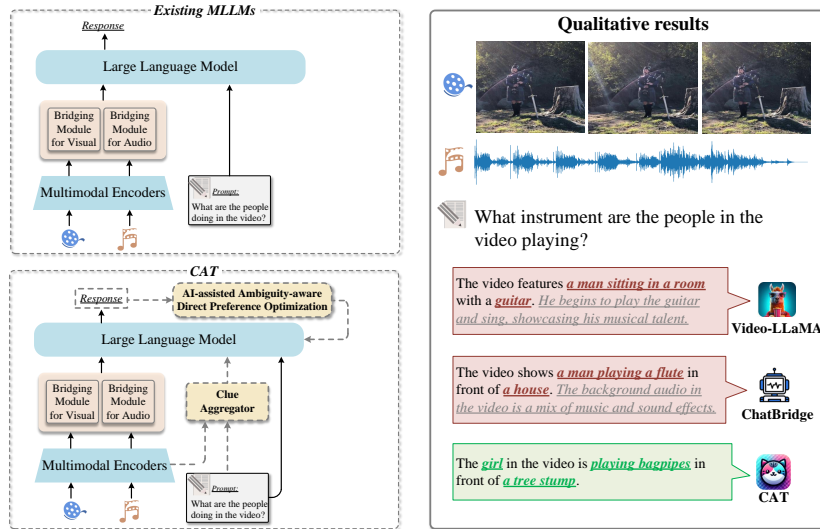


Fig. 1: Comparison between existing MLLMs and CAT. **Red** words for incorrect response, **Green** words for correct response, and **Gray** words for useless response. **Left:** Most of the existing MLLMs straightforwardly bridge multimodal to large language models. Instead, CAT builds on this foundation by designing the clue aggregator itself to learn more detailed knowledge related to the question. Moreover, CAT constrains itself to learn a sharper response through AI-assisted ambiguity-aware direct preference optimization. **Right:** In comparison with audio-visual-language models Video-LLaMA [64] and ChatBridge [67], our method accurately recognizes the answers to questions with the most streamlined responses.

datasets empowers Large Language Models (LLMs) for multimodal understanding. Despite MLLMs’ strong causal reasoning achieving impressive results in common-sense answering [53, 63], academic answering [16], etc., predicting questions in dynamic audio-visual scenarios remains challenging. This is due to the difficulty of aligning LLMs with cross-domain data during training on large-scale multimodal corpora. It leads to particular ambiguity when describing specific objects in dynamic audio-visual scenarios. These interferences not only cause the model to make random guesses but also seriously affect the inference of detailed answers in Audio-Visual Question Answering (AVQA) [1, 26, 59, 62] tasks.

A series of bridging modules [12, 20, 31, 32, 39] for MLLMs have been designed. The simplest methods [15, 18, 28, 33] use projectors to directly align text with other modalities, but partially limit the ability to capture fine-grained information. In addition, using a cross-attention mechanism to query the audio-visual context [8, 21, 38, 39] to solve multimodal alignment problems is effective but still occurs that a certain visual object or sound cannot be localized in practice.

In this work, we explore possible factors contributing to the above failures: **1) Audio-visual insufficiently correlates with the question.** As illustrated in the upper left of Fig. 1, most existing MLLMs [28, 64, 67] are mainly designed with multiple branches to handle multiple modalities individually, followed by

concatenating the modality embeddings with prompts as inputs to the LLM. This paradigm fails to allow textual information to interact with audio-visual at a low level, resulting in the network’s inability to focus on details relevant to the question. **2) Alignment of multimodal with text is challenging.** The multimodal-text corpus [67] is difficult to be aligned, making the model generation of responses sometimes ambiguous. This ambiguity generally manifests in ambiguous words of the corresponding audio-visual content, as well as in generating useless text that overly responds. For example, the descriptions illustrated by Video-LLaMA [64] and ChatBridge [67] in Fig. 1 are all relatively ambiguous in describing the question-related video and audio content. These descriptions fail to identify the correct answer for “bagpipes” when asked for audio-visual reasoning. In addition, when the description is accompanied by many useless and controversy-prone words, it is unfavorable to evaluate the closed-ended AVQA task or open-ended AVQA task.

To overcome the two aforementioned issues, we introduce the CAT, enhancing MLLM in three ways: **1) Aggregation of question-related key clues.** As illustrated in the lower left of Fig. 1, besides the global visual and audio information, we design a clue aggregator to dynamically capture question-aware visual and audio hidden features to enrich fine-grained clues. The aggregator receives sufficient low-level textual information to interact with audio-visual, it is capable of improving audio-visual grounding. **2) Mixed audio-visual training strategy.** The training of CAT includes feature alignment using video-text pairs and audio-text pairs, and high-quality instructions to enhance audio-visual awareness. This strategy allows the CAT to be directly involved in real-world scenarios containing both visual and sound. Notably, we collect an audio-visual joint instruction dataset, named AVinstruct, to further empower CAT in AVQA tasks. **3) Retraining MLLMs to mitigate ambiguity.** The DPO proposed by Rafailov et al. [52] has inspired us. Although MLLMs after training are equipped with multimodal understanding, they lack the flexibility to perceive ambiguity. Therefore, we reframe ambiguity elimination as a model preference optimization process and propose an AI-assisted Ambiguity-aware Direct Preference Optimization (ADPO) strategy. Specifically, we refer to ambiguous responses that express the lack of clarity of specific audio-visual objects as negative responses, and then we collect negative responses in the training set and utilize GPT to rewrite them into positive responses. After the multimodal training, we perform ADPO to retrain the model to bias towards the positive response, which is the precise description after the rewrite, and reject the negative response, which is the ambiguous description. Through this learning strategy, CAT can constrain itself to favor non-ambiguity responses. As shown in Fig. 1 on the right, CAT correctly responds to the question and excludes all useless information.

Our main contributions are summarized as follows:

- We introduce a novel audio-visual-language model, dubbed as CAT, that is capable of learning question-related clues and engaging directly in dynamic audio-visual inference. Notably, we collect AVinstruct, an audio-visual joint instruction dataset to ensure the stability of CAT in AVQA tasks.

- As a powerful learning strategy, we propose AI-assisted ambiguity-aware direct preference optimization to overcome the problem that MLLMs tend to ambiguously describe specific audio-visual objects.
- We evaluate CAT on a wide range of multimodal tasks. Extensive experiments demonstrate the superiority of CAT (e.g., outperforms the state-of-the-art on a variety of AVQA tasks, and achieves remarkable results in the evaluation of video-based text generation tasks and zero-shot video question-answering tasks).

2 Related Works

Audio-visual question answering (AVQA). AVQA aims to produce the most accurate linguistic representation of a given video and audio based on the question content, which requires multimodal understanding and reasoning across different semantic levels. The earliest studies [57, 58, 60] emphasize the understanding of the whole video and return a simple word that is relatively correct to the question. Subsequently, the emergence of dynamic audio-visual datasets (e.g., music scenes [26, 59], and 360-degree panoramas [62]) increases the challenge of Question-Answering (QA), which requires the mining of temporal and spatial information in different modalities. For instance, Music-AVQA [26] requires distinguishing how many instruments are present based on the audio content. AVQA [59] requires finding the most plausible of multiple options based on the audio-visual clues from dynamic scenarios.

MLLMs for audio-visual question answering. Extending LLM to other multimodal tasks is a recently emerging field. Despite the ability of MLLMs to combine information between different modalities, performance on downstream tasks, especially AVQA, remains suboptimal. Many works [8, 20, 39] emphasize the design of elegant bridging methods to improve the performance of question answering. The simplest bridging modules [37, 40, 64] use one or more linear layers for feature alignment. Although such methods minimize parameter updates, they still have limitations in exploring fine-grained information. Others design more complex bridging networks to query visual information. For example, Lyu et al. [38] propose an alignment module for harmonizing different representations before entering the LLM. Ma et al. [39] propose to keep the distance between all visual tokens and any linguistic tokens consistent within the LLM. However, a large multimodal corpus is difficult to align during training, neither algorithms with low parameter counts [12] nor complex bridging networks [4] are prone to the problem of failing to accurately depict specific audio-visual events.

Human-preference learning. Reinforcement Learning from Human Feedback (RLHF) [41, 46] is the most classical instance of human preference learning [45, 55], it constructs a reward model to optimize the policy model to favor preference responses. Dai et al. [8] have demonstrated that such preference learning can enhance LLM to generate more accurate information. Recently, Rafailov et al. [52] propose a Direct Preference Optimization (DPO) strategy that learns

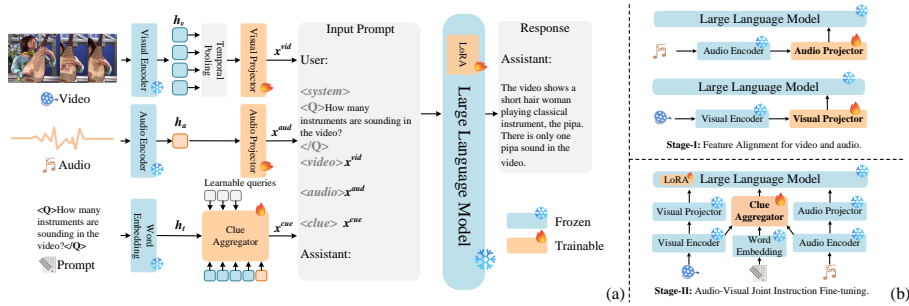


Fig. 2: Illustration of the proposed CAT and its training strategy. (a) Overview of CAT. CAT first extracts overall audio-visual knowledge from video and audio and transforms them into visual tokens x^{vid} and audio tokens x^{aud} . We input question tagged with $\langle Q \rangle \langle /Q \rangle$ in the prompt into the clue aggregator, aiming to aggregate question-aware audio-visual hidden features and yield clue tokens x^{clue} . Finally, we merge multimodal tokens and language and feed into the frozen large language model with LoRA [17] to output the response. (b) The training paradigm of CAT involves pre-alignment of the audio-visual projectors and instruction tuning on the entire model.

preferences directly by bypassing learning reward models, this simple yet efficient approach inspires us to solve the ambiguity problem.

3 Our Approach

In this section, we first present in detail the proposed CAT (Sec. 3.1 and 3.2). As shown on the left in Fig. 2, CAT consists of three branches that draw on visual knowledge, audio knowledge, and question-related clues to feed into the LLM, respectively. Second, we present the multimodal training strategy for CAT as shown on the right in Fig. 2, where a high-quality audio-visual joint instruction dataset is collected for further fine-tuning (Sec. 3.3). Lastly, we introduce the AI-assisted Ambiguity-aware Direct Preference Optimization (ADPO) strategy, which reinforces CAT to favor non-ambiguous descriptions (Sec. 3.4).

3.1 Multimodal Inputs

ImageBind [13] with a single joint embedding space demonstrates strong modality learning ability. Therefore, we leverage the frozen ImageBind as a universal encoder for all modalities. Benefiting from the robust pre-alignment of both visual and audio to text in the encoder, CAT can easily transfer cross-domain knowledge. Given a video \mathbf{V} and an audio \mathbf{A} , the encoded multimodal hidden features can be obtained by:

$$h_v = \text{ImageBind}(\mathbf{V}), \quad h_a = \text{ImageBind}(\mathbf{A}), \quad (1)$$

where $h_v \in \mathbb{R}^{T \times d_h}$, $h_a \in \mathbb{R}^{1 \times d_h}$ are video and audio features, respectively. T denotes the length of the given video and d_h is the specific dimension. Notably,

we temporally compress the frame-level features h_v to address the computational burden. Further, we employ two linear projection layers to align the hidden dimensions of the inputs h_v, h_a to obtain visual tokens x^{vid} and audio tokens x^{aud} , respectively.

3.2 Aggregating Key Clues

An important component of MLLMs is the design of efficient bridging modules. The simple approach is to apply linear layers to incorporate the full visual and audio. The advantage is that it does not require any reduction of spatio-temporal information. However, such bridge approaches sometimes make LLM fail to reflect certain content and consistently generate a chunk of description about the video [28, 64]. Therefore, we devise a Clue Aggregator (CA) that enriches LLM-acquired knowledge by mining multimodal clues related to the question. CA is divided into two steps, as illustrated in Fig. 3.

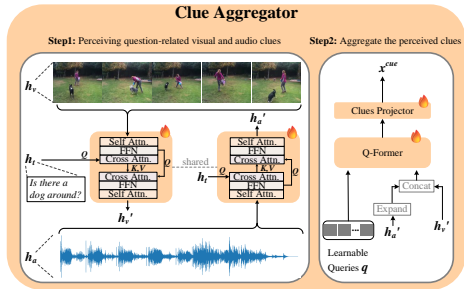


Fig. 3: Illustration of clue aggregator in a simple example.

Step1: perceiving question-related visual and audio clues.

To extract more sufficient details from the input video and audio, we devise a perceiver, which consists of multiple different transformer layers. The perceiver can be viewed as a tiny transformer network that arranges the self-attention (SA), cross-attention (XA), and feed-forward network (FFN) in forward and reverse order. The forward order block \mathcal{B}_1 is to perform attention-based question-aware localization. The reverse order block \mathcal{B}_2 is to consolidate the original audio-visual

representation at the attentional level. Specifically, given tokens representing “*Is there a dog around?*” and the corresponding text embedding h_t :

$$\mathcal{B}_1(h_t; X) = \text{XA}(h_t, \text{FFN}(\text{SA}(X))), \quad (2)$$

$$\mathcal{B}_2(X; \mathcal{B}_1(h_t; X)) = \text{SA}(\text{FFN}(\text{XA}(\text{FFN}(\text{SA}(X)), \mathcal{B}_1(h_t, X))). \quad (3)$$

In $\text{XA}(\cdot, \cdot)$, the former represents the query and the latter represents the key and value. Based on Eqs. 2 and 3, we use $h'_v = \mathcal{B}_2(h_v; \mathcal{B}_1(h_t; h_v))$ and $h'_a = \mathcal{B}_2(h_a; \mathcal{B}_1(h_t; h_a))$ to obtain h'_v, h'_a , the question-aware visual and audio features, respectively. Notably, the perceiver for visual and the perceiver for audio implement shared parameters to learn potential associations.

Step2: aggregating the perceived clues. A three-minute video undergoes a frozen encoder to get frame-level features of about 400 to 500 in length, leading to the fact that general machines cannot bear the burden of feeding all the frame-level features into the LLM. Thanks to the Q-former architecture proposed by Li et al. [27], which reduces the computational cost of end-to-end training of

multimodal-language models. Specifically, we customize a learnable query vector q in length K to further extract useful information from the input question-aware features. Notably, we expand the h'_a in the time dimension to be consistent with h'_v and concatenate them. In practice, we set $K = 48$ to aggregate all question-aware features and obtain the clue tokens x^{cue} via a clues projector to align the dimension with the LLM.

3.3 Multimodal Training Strategy

We follow the previous works [35] to promote comprehension over video and audio. As shown in Fig. 2 on the right, we pre-training CAT on the video-level and audio-level tasks in stage-I. In stage-II, we fine-tune CAT on high-quality audio-visual-level instructions.

Stage-I: feature alignment for video and audio projectors. First, we employ the video-text Webvid 2.5M dataset [3] to train the visual projector (audio projector not involved in training). Second, we employ the audio-text WavCap dataset [42] to train the audio projector (visual projector not involved in training). During the above two training periods, we freeze the LLM and the encoder to align the semantic information for vision and audio, respectively.

Stage-II: audio-visual joint instruction tuning. To equip CAT with the ability to reason jointly based on visual and auditory, we collect an audio-visual joint instruction dataset, named AVinstruct, which emphasizes co-learning dynamic audio-visual pairs to solve diverse AVQA tasks. Specifically, we collect a large number of raw videos containing audio information from YouTube and VGGSound [5] as well as QA pairs from the training set in closed-ended AVQA tasks (i.e. music scene [26] and real-world scene [59]). Then different subtitle generators (i.e. BLIP2 [27] and Whisper [51]) are utilized to obtain video descriptions. Finally, we use GPT to generalize from human-written examples to synthesize question-guided audio-visual descriptions. We integrate all compositions after stage-I training and freeze the visual projector and audio projector, only fine-tuning clue aggregator and combined LoRA parameters on 100k video instruction [40] and AVinstruct. Notably, to highlight the question for applying in the clue aggregator, we reconstruct the form of the input prompt by adding two simple tokens $\langle Q \rangle$, $\langle /Q \rangle$, where $\langle Q \rangle$ denotes the start position of the question and $\langle /Q \rangle$ denotes the end position of the question. During the instruction fine-tuning phase, we use predefined prompts based on the following template:

USER :< *system* >< Q >< $/Q$ >< *video* >< *audio* >< *clues* > Assistant :

In this prompt, $\langle \textit{system} \rangle$ denotes the official guidance message. $\langle \textit{video} \rangle$, $\langle \textit{audio} \rangle$, and $\langle \textit{clues} \rangle$ are associated with visual tokens, audio tokens, and question-related clue tokens, respectively.

3.4 AI-assisted Ambiguity-aware Direct Preference Optimization

At first, we attempt to design various prompts to allow MLLMs to generate the most concise descriptions. However, such a no-learning approach is challenging

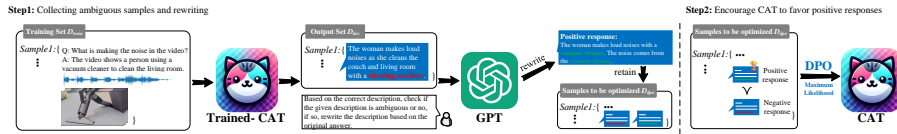


Fig. 4: Trained-CAT denotes CAT after feature alignment and instruction tuning. Our proposed ADPO strategy involves two steps. First, we collect the negative response generated by trained-CAT and correct it by GPT to obtain a positive response based on the original answer. Second, we perform ADPO training to skew CAT toward positive responses and reject negative responses.

to achieve the goal of accurate responses in a variety of dynamic audio-visual scenarios. We find that MLLMs are prone to ambiguity when describing specific audio-visual objects. For example, an ambiguous response is illustrated by trained-CAT in Fig. 4, where trained-CAT denotes CAT after feature alignment and instruction tuning. ADPO is designed to allow MLLMs to review their mistakes and relearn with a new objective to better expression. Specifically, ADPO is divided into two steps to optimize CAT:

Collecting ambiguous samples and rewriting. As illustrated on the left in Fig. 4, assume that a given training set D_{train} and the corresponding trained-CAT output set D_{des} . We first provide GPT with a detailed prompt template to review each output in D_{des} for ambiguity. The rule of the review is to identify output ambiguities, i.e., responses that differ significantly from the original answer, which is referred to as a negative response. We let GPT correct this negative response, with the principle of generating a positive response that clearly describes the visual objects or sound without significant modification. Repeat the above steps to get the sample set D_{dpo} to be optimized, where $D_{dpo} = \{x^{vid}, x^{aud}, x^{cue}, x^{txt}, y_{pos}, y_{neg}\}$. x denotes the input modal token and x^{txt} denotes the text tokens, y_{pos} , y_{neg} denote the positive and negative responses, respectively.

Encouraging CAT to favor positive responses. As illustrated on the right in Fig. 4, ADPO directly optimizes the low-rank adaptation matrix parameters [17] in the policy model. We assume CAT as a policy model f^{pol} , given a reference model f^{ref} , a positive response y_{pos} , and a negative response y_{neg} . The reference model f^{ref} is a deep copy of f^{pol} and during training we encourage f^{pol} to favor y_{pos} , while f^{ref} does not update the weights. Specifically, we define DPO loss \mathcal{L}_{DPO} as:

$$\mathcal{L}_{DPO}(f^{pol}; f^{ref}) = -E_{(x^{vid}, x^{aud}, x^{cue}, x^{txt}, y_{pos}, y_{neg}) \sim D_{dpo}} \left[\log \sigma \left(\beta \log \frac{f^{pol}(y_{pos} | [x^{vid} : x^{aud} : x^{cue} : x^{txt}])}{f^{ref}(y_{pos} | [x^{vid} : x^{aud} : x^{cue} : x^{txt}])} \right) - \beta \log \frac{f^{pol}(y_{neg} | [x^{vid} : x^{aud} : x^{cue} : x^{txt}])}{f^{ref}(y_{neg} | [x^{vid} : x^{aud} : x^{cue} : x^{txt}])} \right), \quad (4)$$

[:] denotes the concatenation. σ is the non-linear function, Sigmoid, and β is a hyperparameter. This objective function aims to directly optimize the

model that favors the positive response y_{pos} and rejects the negative response y_{neg} . Assuming that $\beta = 1$, we simplify the output \hat{r} of the model based on Eq. 4 as:

$$\hat{r}(x^{vid}, x^{aud}, x^{cue}, x^{txt}, y) = \log \frac{f^{pol}(y | [x^{vid} : x^{aud} : x^{cue} : x^{txt}])}{f^{ref}(y | [x^{vid} : x^{aud} : x^{cue} : x^{txt}])}. \quad (5)$$

The optimal goal is to make the variance of $\hat{r}(x^{vid}, x^{aud}, x^{cue}, x^{txt}, y_{pos}) - \hat{r}(x^{vid}, x^{aud}, x^{cue}, x^{txt}, y_{neg})$ larger and larger as a way to steer the model to generate descriptions devoid of ambiguous. However, when the difference between positive and negative responses is small, DPO loss supervision alone is not an effective facilitator. Therefore, we introduce an additional objective \mathcal{L}_{SFT} , which is similar to the process of supervised fine-tuning, to supervise the model to stabilize the bias towards the positive response. Specifically, we define \mathcal{L}_{SFT} as:

$$\mathcal{L}_{SFT} = - \sum \log P(y_{pos} | [x^{video} : x^{audio} : x^{clues} : x^{txt}] ; f^{pol}), \quad (6)$$

$$\{x^{vid}, x^{aud}, x^{cue}, x^{txt}, y_{pos}\} \sim D_{dpo},$$

this loss function can be interpreted as using the positive output probability in the policy model to aid in optimization. During ADPO training, we sum up \mathcal{L}_{DPO} and \mathcal{L}_{SFT} with a bias λ to achieve direct preference optimization:

$$\mathcal{L} = \mathcal{L}_{DPO} + \lambda \mathcal{L}_{SFT}, \quad (7)$$

where we set $\lambda = 0.1$ by default.

4 Experiments

4.1 Datasets

Video-based text generation tasks. We evaluate the proposed CAT on the video understanding benchmarks proposed by Video-ChatGPT [40]. Specifically, Video-ChatGPT proposes five metrics: *Correctness of Information*, *Consistency*, *Detail Orientation*, *Contextual Understanding*, and *Temporal Understanding*, which test the ability of MLLMs to describe videos.

Zero-shot on video question answering tasks. To evaluate whether CAT has the basic ability to communicate regarding video, we conduct zero-shot tests on MSRVTT-QA [58] and ActiviytNet-QA [60]. MSRVTT-QA and ActiviytNet-QA consist of 10k and 5.8k videos containing audio information, respectively, where the QA pairs are mostly questions about daily life.

Closed-ended AVQA tasks. We categorize the Music-AVQA [26] and AVQA [59] datasets as closed-ended AVQA tasks. These datasets consist of up to 42 candidate answers that require the selection of the most appropriate one based on visual and auditory content.

Open-ended AVQA tasks. We select audio-visual dialogue (AVSD [1]), and audio-visual captioning (VALOR [6]) tasks for the evaluation of open-ended

Table 1: GPT-based evaluation [40] for video-based text generation and zero-shot video question answering. For clarity, five scores are reported ("Cr.": Correctness of Information, "Cs.": Consistency, "De.": Detail Orientation, "Ct": Contextual Understanding, "Te.": Temporal Understanding). Score indicates the confidence.

Method	LLM Size	Video-based text generation					Zero-shot on MSRVT-TQA		Zero-shot on ActivityNet-QA	
		Cr.	De.	Ct.	Te.	Cs.	Acc.	Score	Acc.	Score
LLaMA-VID [31]	13B	61.4	61.0	72.0	51.6	52.6	58.9	3.3	47.5	3.3
Video-LLaMA [64]	7B	39.2	43.6	43.2	36.4	35.8	29.6	1.8	12.4	1.1
LLaMA-Adapter [65]	7B	40.6	46.4	46.0	39.6	43.0	43.8	2.7	34.2	2.7
VideoChat [28]	7B	44.6	50.0	50.6	38.8	44.8	45.0	2.5	26.5	2.2
Video-ChatGPT [40]	7B	48.0	50.4	52.4	39.6	47.4	49.3	2.8	35.2	2.7
VISTA-LLaMA [39]	7B	48.8	52.8	63.6	45.2	46.2	60.5	3.3	48.3	3.3
VideoChat2 [29]	7B	60.4	57.8	70.2	53.2	56.2	54.1	3.3	49.1	3.3
Chat-UniVi [20]	7B	57.8	58.2	69.2	57.8	56.2	54.6	3.1	45.8	3.2
LLaMA-VID [31]	7B	59.2	60.0	70.6	49.2	50.2	57.7	3.2	47.1	3.3
CAT (Ours)	7B	61.6	62.0	<u>69.8</u>	<u>56.2</u>	57.8	62.1	3.5	50.2	3.5

AVQA. These tasks require precise language to interpret, correlate, and reason about cross-modal information. We evaluate the zero-shot ability of CAT on these datasets.

4.2 Experimental Setup

Evaluation metrics. For video-based text generation and zero-shot on video question-answering tasks, we follow the evaluation pipeline proposed by VideoChatGPT [40], which is based on GPT-3.5 to evaluate predictive descriptions against correct descriptions, with a score from 0 to 5 indicating accuracy. In this paper, we standardize 0 to 5 as 0 to 100 to align common accuracy rubrics. For closed-ended AVQA tasks, we report the accuracy of the correct sample. For open-ended AVQA tasks, we report the CIDEr [56] that specializes in evaluating visually descriptive tasks.

Architecture. We use frozen ImageBind [13] and LLaMA2-7B [55] as audio-visual encoders and LLM, respectively. The size of modality embeddings for each modality are $\mathbb{R}^{T \times 1024}$. The outputs x^{vid} , x^{aud} , and x^{cue} are $\mathbb{R}^{1 \times 4096}$, $\mathbb{R}^{1 \times 4096}$, and $\mathbb{R}^{48 \times 4096}$, respectively.

Training details. We complete feature alignment training, instruction tuning, and ADPO training with 1 NVIDIA A100 GPU. In detail, For the feature alignment training and instruction tuning, we use the AdamW optimizer with a cosine learning rate decay and a warm-up period. When LoRA is added, we set $r = 64$ and $\alpha = 128$ for the LoRA parameters, and the total batch size is set to 128 for training 1 epoch with a learning rate of $2e^{-5}$. For ADPO training, we only select the training set in the instruction tuning phase for optimization. We set $r = 64$ and $\alpha = 16$ for the LoRA parameters, and the total batch size is set to 1 for training 1 epoch with a learning rate of $4e^{-6}$. The hyperparameter β is set to 0.1.

Table 2: Comparison with non-LLMs-based methods on fine-tuned Music-AVQA dataset.

Method	Language Model	Trainable Params (M)	Audio avg.	Visual avg.	Audio-Visual avg.	Overall avg.
FCNLSTM [11]	CLIP [50]	9.79	68.9	56.2	60.4	60.8
GRU [2]	CLIP	-	68.3	67.0	63.0	65.0
HCAtn [36]	CLIP	-	64.9	65.3	60.3	62.5
MCAN [61]	CLIP	56.0	70.6	71.8	61.5	65.8
PSAC [30]	CLIP	-	72.0	69.4	63.6	66.6
HME [10]	CLIP	-	69.9	68.8	64.8	66.7
HCRN [24]	CLIP	-	63.7	65.2	49.8	56.3
AVSD [54]	CLIP	8.35	68.8	70.3	65.4	67.3
Panp-AVQA [62]	CLIP	-	72.1	73.2	67.0	69.5
AVST [26]	CLIP	18.48	73.9	74.4	69.5	71.6
PSTP-Net [25]	BERT [9]	4.297	70.9	77.3	72.6	73.5
LAVISH [34]	CLIP	21.09	77.1	77.3	77.0	-
CAD [43]	GLoVe [48]	-	78.1	79.7	76.9	78.2
VALOR [6]	BERT	-	-	-	-	78.9
VAST [7]	BERT	-	-	-	-	80.7
CAT-7B (Ours)	LLaMA2	5.813	84.9	86.1	83.2	84.3

4.3 Comparison to State-of-the-Art

Comparison on video-based text generation tasks. We follow the benchmark proposed by Video-ChatGPT [40] to evaluate CAT. On the left of Table 1, we show that CAT achieves the state-of-the-art in terms of correctness of description information (Cr.), detailed description of the problem (De.), and coherence of the description (Cs.). In addition, CAT does not parse longer visual tokens, it still achieves competitive results when comparing time-related descriptions (e.g., contextual and temporal understanding).

Comparison on zero-shot video question answering tasks. On the right of Table 1, we show the zero-shot video question answering performance of CAT on several open-ended datasets. While recent MLLMs designed with bridging modules have produced substantial results, CAT is way ahead of them in recognizing accurate answers. We consistently outperform the state-of-the-art on the MSRVT-T-QA [58] and ActivityNet-QA [60] benchmarks.

Comparison on closed-ended AVQA tasks. We choose to evaluate Music-AVQA [26] to demonstrate that CAT is capable of perceiving specific audio-visual objects to answer questions. In Table 2, under full supervision of the training set, CAT accurately retrieves specific objects in dynamic audio-visual scenarios and comprehensively outperforms all non-LLMs-based models. Thanks to LoRA [17], we can improve the evaluation quality with the help of the world knowledge inside LLaMA2 [55] and only 5.8M parameters are trainable. Also, we examine the ability of CAT for zero-shot on Music-AVQA [26]. In Table 3, we show the results comparing LLMs-based models. Notably, for a fair comparison, we remove the LoRA parameters of CAT that are fine-tuned on AVinstruct, which is derived from the Music-AVQA [26] and AVQA [59] training sets. Even though our model does not draw on the larger-scale LLM, it still achieves a small advantage over ChatBridge [67] with a 13B LLM size. Furthermore, we show the

Table 3: Comparison with LLMs-based methods on zero-shot Music-AVQA dataset.

Method	zero-shot	Acc.
OneLLM-7B [14]	✓	43.0
ChatBridge-13B [67]	✓	47.6
CAT-7B (Ours)	✓	48.6

Table 4: Evaluation on fine-tuned AVQA dataset.

Method	Acc.
HGA + HAVF [19]	87.7
HCRN + HAVF [59]	89.0
PSTP-Net [25]	90.2
CAT-7B (Ours)	92.0

Table 5: Evaluation on zero-shot open-ended AVQA datasets.

Method	zero-shot	AVSD CIDEr	VALOR CIDEr
VALOR [6]	×	-	61.5
VAST [7]	×	-	62.2
FA+HRED [44]	×	84.3	-
MTN [23]	×	98.5	-
COST [49]	×	108.5	-
OneLLM-7B [14]	✓	74.5	29.2
ChatBridge-13B [67]	✓	75.4	24.7
CAT-7B (Ours)	✓	79.0	32.4

performance of CAT in multiple-choice scenarios [59] in Table 4, CAT continues outstanding.

Comparison on open-ended AVQA tasks. Open-ended AVQA tasks require responses to daily events based on audio-visual content. We conduct a comparative analysis of our model with multimodal-based LLMs: OneLLM [14], and ChatBridge [67]. The tests are divided into zero-shot and full supervision, and we only choose to evaluate the zero-shot to demonstrate the level of practical application of the CAT. In Table 5, CAT surpasses OneLLM and ChatBridge on zero-shot complex text reasoning tasks AVSD [1] and VALOR [6]. Moreover, our approach is close to the results of FA [44] using fully supervised training, further demonstrating the strong multimodal understanding of CAT.

4.4 Ablations and Analyses

In this subsection, we explore the effects of each component of CAT. We test the effect of input modal tokens on video-based text generation [40]. In addition, to avoid the huge expense associated with the GPT3.5 evaluation, we test each component of the clue aggregator on the test set of AVinstruct, as well as the generalizability of ADPO on AVSD [1].

Does question-related clues matter? Visual tokens and audio tokens already have the ability to distinguish video content. To explore whether the added question-related clues are meaningful, we conduct ablation experiments on the video-based text generation task in Table 6. We first test visual knowledge separately from question-related clues knowledge on video comprehension. Observations reveal that the results achieved when inputting question-related clues alone have been able to outperform the joint input of audio-visual knowledge. Further, we add clues knowledge to the original audio-visual knowledge, and the results demonstrate that the fusion of the three can fully satisfy the information required for LLM reasoning.

Table 6: Ablation experiment about input modals on video-based text generation. **Table 7:** Ablation experiments about clue aggregator on audio-visual joint instruction dataset. **Table 8:** Ablation experiment about ADPO on AVSD [1].

Input-modal			Cr.	De.	Te.
x^{vid}	x^{aud}	x^{cue}			
✓			38.6	40.4	34.0
✓	✓		40.6	44.8	35.2
		✓	58.6	57.4	55.8
✓	✓	✓	61.6	59.0	56.2

Method	B@4	M	R
$\mathcal{B}_1 + \mathcal{B}_2$	30.8	58.7	68.5
$\mathcal{B}_1 + \mathcal{B}_2$ w/o h_v	10.0	37.1	49.9
$\mathcal{B}_1 + \mathcal{B}_2$ w/o h_a	18.8	49.3	61.8
only \mathcal{B}_1	25.9	54.7	64.2

Model	ADPO	AVSD		
		B@4	M	C
Video-LLaMA	×	18.4	40.1	63.7
Video-LLaMA	✓	22.2	48.2	69.2
CAT (Ours)	×	28.9	56.2	74.8
CAT (Ours)	✓	34.2	59.8	79.0

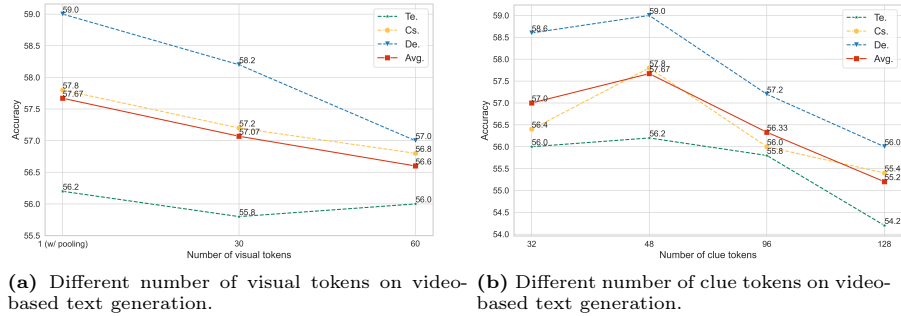


Fig. 5: The impacts of input modal tokens. Avg. represents the average accuracy of temporal (Te.), consistency (Cs.), and detail (De.).

Ablation on clue aggregator. We categorize a portion of the collected audio-visual joint instruction data into a test set and present the ablation results in Table 7. The evaluation indicators B, M, and C denote BLEU-4 [47], METEOR [22], and ROUGE-L, respectively. We split \mathcal{B}_1 and \mathcal{B}_2 and query subjects h_v and h_a to explore the impact of CA internals on reasoning about dynamic audio-visual pairs. We find that visual attention influences context much more than auditory. This is because vision has more information that can be mined and is more relevant to the question. Next, we reduce parameters to investigate whether the same performance can be achieved with a single block. However, since the question acts as a query, without \mathcal{B}_2 to recover the original modal length does not provide an advantage over two blocks.

Number of input modal tokens. We evaluate the effect of the number of input modal tokens on video understanding [40]. As shown in Fig. 5, we first study the effect of the number of visual tokens on a time-related description task. Indeed, as the number of visual tokens increases, the accuracy in the evaluation of the various descriptions decreases, even if the change is not large, but it also confirms that the language model may not be able to reason based on longer visual tokens. Next, we examine the impact of the number of clue tokens set by the Q-former [27]. Clue aggregator uses Q-Former to transform frame-level features into specific tokens, and we find that it works best when the length of queries K is set to 48.

Impact of ADPO. As the paper explains, supervised fine-tuning alone still underperforms in audio-visual scenarios. We test the generalizability of ADPO

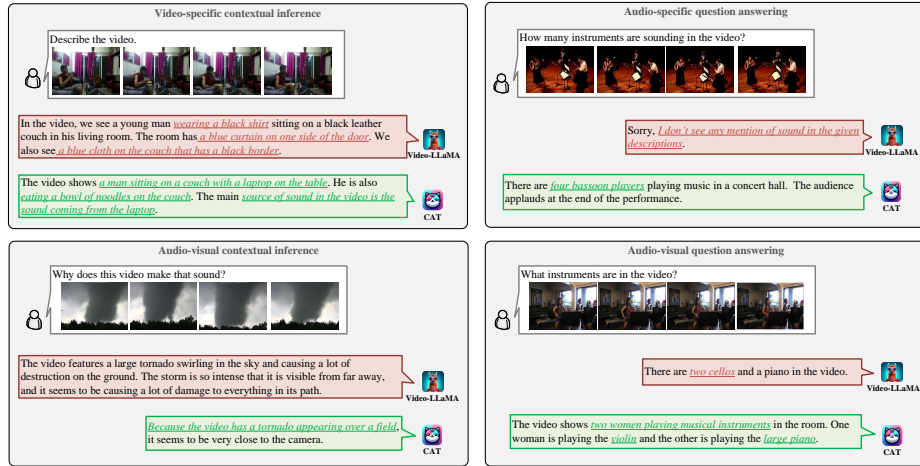


Fig. 6: Qualitative results of the video-specific contextual inference, audio-visual contextual inference, audio-specific question answering, and audio-visual question answering, respectively.

on Video-LLaMA [64] and ADPO’s effect on AVSD [1] in Table 8 to demonstrate the superiority of this learning strategy. C denotes the CiDER [56]. We find that for different MLLMs, ADPO does improve the descriptive power of the model without much learning cost, it brings different levels of enhancement to both Video-LLaMA and CAT.

4.5 Qualitative Analysis

In Fig. 6, we analyze the qualitative results with Video-LLaMA [64], which is also based on the audio-visual-language model. In the example of video-specific contextual inference, we show that CAT has a sharp perception of complex indoor scenes. Our descriptions can accurately represent what the person doing and what background sounds are in the video. Video-LLaMA is biased towards describing scene information and incorrect character information. In the specific audio question answering example, Video-LLaMA lost the ability to answer “How many instruments are sounding in the video?”. In contrast, CAT can accurately answer the quantity question. In the audio-visual question answering example, Video-LLaMA answers incorrectly due to failed audio-visual grounding, while our CAT answers correctly due to its strong ability to capture specific objects in audio-visual scenarios.

5 Conclusion

In this work, we introduce CAT to enhance LLMs’ multimodal understanding in dynamic audio-visual scenarios. We propose a clue aggregator to aggregate the question-related clues, which enriches the knowledge required by LLM for detailed reasoning. We mix datasets containing audio and video to empower LLM

with multimodal understanding. To more consistently infer audio-visual scenarios, we collect an audio-visual joint instruction dataset to further fine-tune the CAT. Moreover, we propose AI-assisted ambiguity-aware direct preference optimization, a strategy specialized in retraining the model to eliminate ambiguous descriptions for more accurate responses to specific audio-visual objects. CAT has achieved comparable results when applied to a variety of complex audio-visual scenarios. We consider further modal extensions in the future for more comprehensive realistic applications.

References

1. AlAmri, H., Cartillier, V., Das, A., Wang, J., Cherian, A., Essa, I., Batra, D., Marks, T.K., Hori, C., Anderson, P., Lee, S., Parikh, D.: Audio visual scene-aware dialog. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR. pp. 7558–7567 (2019) [2](#), [9](#), [12](#), [13](#), [14](#)
2. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: VQA: visual question answering. In: IEEE International Conference on Computer Vision, ICCV. pp. 2425–2433. IEEE Computer Society (2015) [11](#)
3. Bain, M., Nagrani, A., Varol, G., Zisserman, A.: Frozen in time: A joint video and image encoder for end-to-end retrieval. In: IEEE/CVF International Conference on Computer Vision, ICCV. pp. 1708–1718. IEEE (2021) [7](#)
4. Chen, G., Shen, L., Shao, R., Deng, X., Nie, L.: LION : Empowering multimodal large language model with dual-level visual knowledge. CoRR [abs/2311.11860](#) (2023) [4](#)
5. Chen, H., Xie, W., Vedaldi, A., Zisserman, A.: Vggsound: A large-scale audio-visual dataset. In: IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP. pp. 721–725. IEEE (2020) [7](#)
6. Chen, S., He, X., Guo, L., Zhu, X., Wang, W., Tang, J., Liu, J.: VALOR: vision-audio-language omni-perception pretraining model and dataset. CoRR [abs/2304.08345](#) (2023) [9](#), [11](#), [12](#)
7. Chen, S., Li, H., Wang, Q., Zhao, Z., Sun, M., Zhu, X., Liu, J.: VAST: A vision-audio-subtitle-text omni-modality foundation model and dataset. CoRR [abs/2305.18500](#) (2023) [11](#), [12](#)
8. Dai, W., Li, J., Li, D., Tiong, A.M.H., Zhao, J., Wang, W., Li, B., Fung, P., Hoi, S.C.H.: Instructblip: Towards general-purpose vision-language models with instruction tuning. CoRR [abs/2305.06500](#) (2023) [2](#), [4](#)
9. Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. In: Burstein, J., Doran, C., Solorio, T. (eds.) Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT. pp. 4171–4186. Association for Computational Linguistics (2019) [11](#)
10. Fan, C., Zhang, X., Zhang, S., Wang, W., Zhang, C., Huang, H.: Heterogeneous memory enhanced multimodal attention model for video question answering. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR. pp. 1999–2007. Computer Vision Foundation / IEEE (2019) [11](#)
11. Fayek, H.M., Johnson, J.: Temporal reasoning via audio question answering. IEEE ACM Trans. Audio Speech Lang. Process. **28**, 2283–2294 (2020) [11](#)

12. Gao, P., Han, J., Zhang, R., Lin, Z., Geng, S., Zhou, A., Zhang, W., Lu, P., He, C., Yue, X., Li, H., Qiao, Y.: Llama-adapter V2: parameter-efficient visual instruction model. CoRR **abs/2304.15010** (2023) [2](#), [4](#)
13. Girdhar, R., El-Nouby, A., Liu, Z., Singh, M., Alwala, K.V., Joulin, A., Misra, I.: Imagebind one embedding space to bind them all. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR. pp. 15180–15190 (2023) [5](#), [10](#)
14. Han, J., Gong, K., Zhang, Y., Wang, J., Zhang, K., Lin, D., Qiao, Y., Gao, P., Yue, X.: Onellm: One framework to align all modalities with language. CoRR **abs/2312.03700** (2023) [1](#), [12](#)
15. Han, J., Zhang, R., Shao, W., Gao, P., Xu, P., Xiao, H., Zhang, K., Liu, C., Wen, S., Guo, Z., Lu, X., Ren, S., Wen, Y., Chen, X., Yue, X., Li, H., Qiao, Y.: Imagebind-llm: Multi-modality instruction tuning. CoRR **abs/2309.03905** (2023) [2](#)
16. Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., Steinhardt, J.: Measuring massive multitask language understanding. In: ICLR (2021) [2](#)
17. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. In: The Tenth International Conference on Learning Representations, ICLR (2022) [5](#), [8](#), [11](#)
18. Huang, B., Wang, X., Chen, H., Song, Z., Zhu, W.: Vtimellm: Empower LLM to grasp video moments. CoRR **abs/2311.18445** (2023) [2](#)
19. Jiang, P., Han, Y.: Reasoning with heterogeneous graph alignment for video question answering. In: AACL. pp. 11109–11116. AACL Press (2020) [12](#)
20. Jin, P., Takanobu, R., Zhang, C., Cao, X., Yuan, L.: Chat-univi: Unified visual representation empowers large language models with image and video understanding. CoRR **abs/2311.08046** (2023) [2](#), [4](#), [10](#)
21. Korbar, B., Xian, Y., Tonioni, A., Zisserman, A., Tombari, F.: Text-conditioned resampler for long form video understanding. CoRR **abs/2312.11897** (2023) [2](#)
22. Lavie, A., Agarwal, A.: METEOR: an automatic metric for MT evaluation with high levels of correlation with human judgments. In: Callison-Burch, C., Koehn, P., Fordyce, C.S., Monz, C. (eds.) Proceedings of the Second Workshop on Statistical Machine Translation. pp. 228–231. Association for Computational Linguistics (2007) [13](#)
23. Le, H., Sahoo, D., Chen, N.F., Hoi, S.C.H.: Multimodal transformer networks for end-to-end video-grounded dialogue systems. In: Korhonen, A., Traum, D.R., Màrquez, L. (eds.) Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL. pp. 5612–5623. Association for Computational Linguistics (2019) [12](#)
24. Le, T.M., Le, V., Venkatesh, S., Tran, T.: Hierarchical conditional relation networks for video question answering. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR. pp. 9969–9978. Computer Vision Foundation / IEEE (2020) [11](#)
25. Li, G., Hou, W., Hu, D.: Progressive spatio-temporal perception for audio-visual question answering. In: El-Saddik, A., Mei, T., Cucchiara, R., Bertini, M., Vallejo, D.P.T., Atrey, P.K., Hossain, M.S. (eds.) ACM MM. pp. 7808–7816. ACM (2023) [11](#), [12](#)
26. Li, G., Wei, Y., Tian, Y., Xu, C., Wen, J.R., Hu, D.: Learning to answer questions in dynamic audio-visual scenarios. In: CVPR. p. 19086–19096 (2022) [2](#), [4](#), [7](#), [9](#), [11](#)
27. Li, J., Li, D., Savarese, S., Hoi, S.C.H.: BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In: Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., Scarlett, J. (eds.) International Conference on Machine Learning, ICML. pp. 19730–19742. PMLR (2023) [6](#), [7](#), [13](#)

28. Li, K., He, Y., Wang, Y., Li, Y., Wang, W., Luo, P., Wang, Y., Wang, L., Qiao, Y.: Videochat: Chat-centric video understanding. CoRR [abs/2305.06355](#) (2023) [2](#), [6](#), [10](#)
29. Li, K., Wang, Y., He, Y., Li, Y., Wang, Y., Liu, Y., Wang, Z., Xu, J., Chen, G., Luo, P., Wang, L., Qiao, Y.: Mvbench: A comprehensive multi-modal video understanding benchmark. CoRR [abs/2311.17005](#) (2023) [10](#)
30. Li, X., Song, J., Gao, L., Liu, X., Huang, W., He, X., Gan, C.: Beyond rnns: Positional self-attention with co-attention for video question answering. In: AAAI. pp. 8658–8665. AAAI Press (2019) [11](#)
31. Li, Y., Wang, C., Jia, J.: Llama-vid: An image is worth 2 tokens in large language models. CoRR [abs/2311.17043](#) (2023) [2](#), [10](#)
32. Li, Z., Yang, B., Liu, Q., Ma, Z., Zhang, S., Yang, J., Sun, Y., Liu, Y., Bai, X.: Monkey: Image resolution and text label are important things for large multi-modal models. CoRR [abs/2311.06607](#) (2023) [2](#)
33. Lin, B., Ye, Y., Zhu, B., Cui, J., Ning, M., Jin, P., Yuan, L.: Video-llava: Learning united visual representation by alignment before projection. CoRR [abs/2311.10122](#) (2023) [2](#)
34. Lin, Y., Sung, Y., Lei, J., Bansal, M., Bertasius, G.: Vision transformers are parameter-efficient audio-visual learners. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR. pp. 2299–2309. IEEE (2023) [11](#)
35. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. CoRR [abs/2304.08485](#) (2023) [1](#), [7](#)
36. Lu, J., Yang, J., Batra, D., Parikh, D.: Hierarchical question-image co-attention for visual question answering. In: Lee, D.D., Sugiyama, M., von Luxburg, U., Guyon, I., Garnett, R. (eds.) Advances in Neural Information Processing Systems NIPS. pp. 289–297 (2016) [11](#)
37. Luo, R., Zhao, Z., Yang, M., Dong, J., Qiu, M., Lu, P., Wang, T., Wei, Z.: Valley: Video assistant with large language model enhanced ability. CoRR [abs/2306.07207](#) (2023) [4](#)
38. Lyu, C., Wu, M., Wang, L., Huang, X., Liu, B., Du, Z., Shi, S., Tu, Z.: Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration. CoRR [abs/2306.09093](#) (2023) [2](#), [4](#)
39. Ma, F., Jin, X., Wang, H., Xian, Y., Feng, J., Yang, Y.: Vista-llama: Reliable video narrator via equal distance to visual tokens. CoRR [abs/2312.08870](#) (2023) [2](#), [4](#), [10](#)
40. Maaz, M., Rasheed, H.A., Khan, S.H., Khan, F.S.: Video-chatgpt: Towards detailed video understanding via large vision and language models. CoRR [abs/2306.05424](#) (2023) [4](#), [7](#), [9](#), [10](#), [11](#), [12](#), [13](#)
41. MacGlashan, J., Ho, M.K., Loftin, R.T., Peng, B., Wang, G., Roberts, D.L., Taylor, M.E., Littman, M.L.: Interactive learning from policy-dependent human feedback. In: Precup, D., Teh, Y.W. (eds.) Proceedings of the 34th International Conference on Machine Learning, ICML. pp. 2285–2294 (2017) [4](#)
42. Mei, X., Meng, C., Liu, H., Kong, Q., Ko, T., Zhao, C., Plumbley, M.D., Zou, Y., Wang, W.: Wavcaps: A chatgpt-assisted weakly-labelled audio captioning dataset for audio-language multimodal research. CoRR [abs/2303.17395](#) (2023) [7](#)
43. Nadeem, A., Hilton, A., Dawes, R., Thomas, G., Mustafa, A.: CAD - contextual multi-modal alignment for dynamic AVQA. CoRR [abs/2310.16754](#) (2023) [11](#)
44. Nguyen, D.T., Sharma, S., Schulz, H., Asri, L.E.: From film to video: Multi-turn question answering with multi-modal context. CoRR [abs/1812.07023](#) (2018) [12](#)
45. OpenAI: GPT-4 technical report. CoRR [abs/2303.08774](#) (2023) [4](#)

46. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C.L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P.F., Leike, J., Lowe, R.: Training language models to follow instructions with human feedback. In: *Advances in Neural Information Processing Systems, NeurIPS (2022)* **1, 4**
47. Papineni, K., Roukos, S., Ward, T., Zhu, W.: Bleu: a method for automatic evaluation of machine translation. In: *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*. pp. 311–318. *ACL (2002)* **13**
48. Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: Moschitti, A., Pang, B., Daelemans, W. (eds.) *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP*. pp. 1532–1543. *ACL (2014)* **11**
49. Pham, H., Le, T.M., Le, V., Phuong, T.M., Tran, T.: Video dialog as conversation about objects living in space-time. In: Avidan, S., Brostow, G.J., Cissé, M., Farinella, G.M., Hassner, T. (eds.) *ECCV*. vol. 13699, pp. 710–726. Springer (2022) **12**
50. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. In: Meila, M., Zhang, T. (eds.) *Proceedings of the International Conference on Machine Learning, ICML*. vol. 139, pp. 8748–8763. PMLR (2021) **11**
51. Radford, A., Kim, J.W., Xu, T., Brockman, G., McLeavey, C., Sutskever, I.: Robust speech recognition via large-scale weak supervision. In: Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., Scarlett, J. (eds.) *International Conference on Machine Learning, ICML*. pp. 28492–28518. PMLR (2023) **7**
52. Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning, C.D., Finn, C.: Direct preference optimization: Your language model is secretly a reward model. *CoRR abs/2305.18290* (2023) **3, 4**
53. Sakaguchi, K., Bras, R.L., Bhagavatula, C., Choi, Y.: Winogrande: An adversarial winograd schema challenge at scale. In: *AAAI*. pp. 8732–8740 (2020) **2**
54. Schwartz, I., Schwing, A.G., Hazan, T.: A simple baseline for audio-visual scene-aware dialog. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. pp. 12548–12558. Computer Vision Foundation / IEEE (2019) **11**
55. Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bakhlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Canton-Ferrer, C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardaş, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P.S., Lachaux, M., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E.M., Subramanian, R., Tan, X.E., Tang, B., Taylor, R., Williams, A., Kuan, J.X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., Scialom, T.: Llama 2: Open foundation and fine-tuned chat models. *CoRR abs/2307.09288* (2023) **4, 10, 11**
56. Vedantam, R., Zitnick, C.L., Parikh, D.: Cider: Consensus-based image description evaluation. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. pp. 4566–4575 (2015) **10, 14**
57. Wang, X., Wu, J., Chen, J., Li, L., Wang, Y., Wang, W.Y.: VateX: A large-scale, high-quality multilingual dataset for video-and-language research. In: *IEEE/CVF*

- International Conference on Computer Vision, ICCV. pp. 4580–4590. IEEE (2019) [4](#)
58. Xu, D., Zhao, Z., Xiao, J., Wu, F., Zhang, H., He, X., Zhuang, Y.: Video question answering via gradually refined attention over appearance and motion. In: ACM MM. pp. 1645–1653 (2017) [4](#), [9](#), [11](#)
 59. Yang, P., Wang, X., Duan, X., Chen, H., Hou, R., Jin, C., Zhu, W.: AVQA: A dataset for audio-visual question answering on videos. In: ACM MM. pp. 3480–3491. ACM (2022) [2](#), [4](#), [7](#), [9](#), [11](#), [12](#)
 60. Yu, Z., Xu, D., Yu, J., Yu, T., Zhao, Z., Zhuang, Y., Tao, D.: Activitynet-qa: A dataset for understanding complex web videos via question answering. In: AAAI. pp. 9127–9134 (2019) [4](#), [9](#), [11](#)
 61. Yu, Z., Yu, J., Cui, Y., Tao, D., Tian, Q.: Deep modular co-attention networks for visual question answering. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR. pp. 6281–6290. Computer Vision Foundation / IEEE (2019) [11](#)
 62. Yun, H., Yu, Y., Yang, W., Lee, K.I., Kim, G.H.: Pano-avqa: Grounded audio-visual question answering on 360° videos. In: CVPR. p. 2031–2041 (2021) [2](#), [4](#), [11](#)
 63. Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., Choi, Y.: Hellaswag: Can a machine really finish your sentence? In: Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL. pp. 4791–4800 (2019) [2](#)
 64. Zhang, H., Li, X., Bing, L.: Video-llama: An instruction-tuned audio-visual language model for video understanding. In: Proceedings of the Empirical Methods in Natural Language Processing, EMNLP. pp. 543–553 (2023) [1](#), [2](#), [3](#), [4](#), [6](#), [10](#), [14](#)
 65. Zhang, R., Han, J., Zhou, A., Hu, X., Yan, S., Lu, P., Li, H., Gao, P., Qiao, Y.: Llama-adapter: Efficient fine-tuning of language models with zero-init attention. CoRR [abs/2303.16199](#) (2023) [10](#)
 66. Zhang, S., Dong, L., Li, X., Zhang, S., Sun, X., Wang, S., Li, J., Hu, R., Zhang, T., Wu, F., Wang, G.: Instruction tuning for large language models: A survey. CoRR [abs/2308.10792](#) (2023) [1](#)
 67. Zhao, Z., Guo, L., Yue, T., Chen, S., Shao, S., Zhu, X., Yuan, Z., Liu, J.: Chat-bridge: Bridging modalities with large language model as a language catalyst. CoRR [abs/2305.16103](#) (2023) [1](#), [2](#), [3](#), [11](#), [12](#)