Understanding Sounds, Missing the Questions: The Challenge of Object Hallucination in Large Audio-Language Models

Chun-Yi Kuan, Wei-Ping Huang, Hung-yi Lee

Graduate Institute of Communication Engineering, National Taiwan University, Taiwan

chunyi.kaun.tw@gmail.com, r11942102@ntu.edu.tw, hungyilee@ntu.edu.tw

Abstract

Large audio-language models (LALMs) enhance traditional large language models by integrating audio perception capabilities, allowing them to tackle audio-related tasks. Previous research has primarily focused on assessing the performance of LALMs across various tasks, yet overlooking their reliability, particularly concerning issues like object hallucination. In our study, we introduce methods to assess the extent of object hallucination of publicly available LALMs. Our findings reveal that LALMs are comparable to specialized audio captioning models in their understanding of audio content, but struggle to answer discriminative questions, specifically those requiring the identification of the presence of particular object sounds within an audio clip. This limitation highlights a critical weakness in current LALMs: their inadequate understanding of discriminative queries. Moreover, we explore the potential of prompt engineering to enhance LALMs' performance on discriminative questions.

Index Terms: Large audio-language models, Object hallucination

1. Introduction

Large audio-language models (LALMs) augment traditional large language models (LLMs) by incorporating audio perception capabilities. These models can accept both audio and text instructions as inputs, facilitating a broader spectrum of audiorelated tasks. Recently, numerous works [1–17] have proposed integrating audio perception modules and LLMs into a single multimodal model, which enables LLMs to process and understand audio inputs and handle various speech and audio tasks.

In the field of evaluation for LALMs, Dynamic-SUPERB [18] is introduced as the first benchmark focusing on speech processing tasks through designed questions and answer options. Following this, AIR-Bench [19] was developed, employing open-ended questions with GPT-4 [20] as the evaluator to assess LALMs' task performance. Both benchmarks primarily evaluate task performance but do not adequately assess the reliability of the content generated by LALMs, particularly concerning object hallucination. The issue of hallucination [21–23] has been raised with the rapid development of LLMs, and numerous works [24–30] in the field of computer vision have observed that large vision-language models (LVLMs) exhibit significant object hallucination in image captioning tasks, generating captions that include objects not present in the images.

Given the lack of discussion on object hallucination in LALMs within speech and audio domains, and the absence of benchmarks specifically measuring object hallucination, this paper introduces discriminative and generative tasks aimed at exploring the object hallucination phenomenon within LALMs.

Discriminative tasks aim to ascertain the presence of a specific object's sound within an audio clip by asking the model questions such as, "Is there a sound of a dog in the audio?". The design of these questions involves sampling objects through both positive and negative methods. Consequently, we can treat discriminative tasks as a binary classification task, calculating metrics such as accuracy, precision, recall, and the F1 score to evaluate performance. On the other hand, generative tasks are designed to guide models in performing audio captioning tasks through instructions, such as "Describe the audio". After obtaining the model's predicted audio captions, nouns are extracted from these captions using NLP tools. These nouns are considered as objects present in the sound, which are then compared with the ground truth labels to ascertain the extent of object hallucination exhibited by the model.

In the results of these two different types of tasks, we found that LALMs suffer from object hallucination and tend to give affirmative answers. In addition, the performance of these LALMs is highly sensitive to prompt design. Interestingly, LALMs rival specialized models in audio captioning tasks, demonstrating their ability to comprehend audio information. However, their performance in discriminative tasks is less satisfactory. Although the models are adept at audio captioning, they struggle with answering discriminative questions and suffer from object hallucination. We observe that LALMs struggle to understand discriminative questions, failing to extract the required information from the given audio. Thus, we propose methods to improve their performance on discriminative tasks. Our contributions are outlined as follows:

- This is the first work to explore object hallucination in large audio-language models (LALMs).
- We observe that LALMs perform well on audio captioning tasks but struggle with answering discriminative questions, and we propose methods for improvement.

2. Evaluation Methods

In Sec 2.1, we introduce the dataset. Sec 2.2 and 2.3 detail evaluation methods and metrics for discriminative and generative tasks, respectively, as further illustrated in Figure 1.

2.1. Dataset

- Audio: AudioCaps [31] is designed for audio captioning tasks. It is distinguished by its emphasis on human-annotated captions for a wide variety of audio clips from the AudioSet [32], which features labels for different sound events and categories identified within these clips. Therefore, each audio clip is accompanied by text captions and multiple labels.
- Speech: CHIME-6 [33] is tailored for advancing ASR sys-



Figure 1: Demonstration of our evaluation pipeline.

tems in noisy conditions. We use it for noisy speech recognition to validate the models' capabilities in ASR under conditions with environmental noise interference. To ensure that the transcriptions feature a sufficient number of nouns, we selected instances in the test set with more than three nouns in the transcription, totaling 489 instances.

2.2. Discriminative Tasks

Methods: Inspired by Polling-based Object Probing Evaluation (POPE) [24], we adopt the similar approach for evaluating LALMs. We formulate the evaluation of object hallucination as a binary classification task that prompts LALMs to output "Yes" or "No". We design five different prompts, "Is there a sound of [object]?", "Does the audio contain the sound of [object]?", 'Have you noticed the sound of [object]?", "Can you hear the sound of [object]?" and "Can you detect the sound of [object]?". In this way, by utilizing different strategies to sample objects that LALMs prone to hallucinate, we can establish a set of questions to poll LALMs. Since the expected answers to these discriminative questions are simply "Yes" or "No", we can easily identify them without complicated parsing rules.

Questions with answers are "Yes" can be directly built from ground truth objects, while questions whose answers are "No" can be built by sampling from negative objects. Hence, we can devise various sampling strategies to validate whether LALMs are prone to hallucinate the specific objects. We consider all the ground truth labels for each audio as positive samples. Inspired by [24], we utilize the following negative sampling strategies:

- Random Sampling: We random sample *k* objects that are not present in the current audio.
- Popular Sampling: We select the top-k most frequent objects across the entire audio captioning dataset that are not present in the current audio.
- Adversarial Sampling: We rank all objects by their frequency
 of co-occurrence with the actual objects in the audio. Subsequently, we select the top-k objects with the highest cooccurrence rates that are not included in the audio.

To ensure a balanced ratio between positive and negative samples during data construction, the parameter k is adjusted to match the number of ground truth labels associated with each audio clip. Hence, we derive 15,110 positive instances directly from the ground truth. For each negative sampling strategy, we sample the equivalent number of negative instances, which is 15,110, to maintain this balance.

Metrics: We adopt accuracy, precision, recall and F1 score as the evaluation metrics. Because the evaluation is aimed at hal-

lucination, both precision and recall are calculated in relation to hallucination questions, where the ground truth answer is "No". In addition, we report the ratio that LALMs answer "Yes" as a reference to analyze the behavior of models.

2.3. Generative Tasks

Methods: We devise five different prompts, "Describe the audio", "What do you hear?", "What can be inferred from the audio?", "This is a sound of?" and "Generate audio caption:", to prompt LALMs to generate captions for given audio clips. Next, we use NLP tools, SpaCy [37], to extract nouns from these captions, identifying these nouns as the objects that the model perceives to be producing sounds in the audio. This approach allows us to compile a list of objects identified by the LLM as present in the audio. On the other hand, Audio-Caps [31] dataset provides corresponding ground truth captions and labels (lists of objects producing sounds in the audio). By applying the same method to extract nouns from these ground truth captions and then combining them with the ground truth labels, we can obtain a comprehensive list of ground truth objects. By employing the same process, we can adapt it to tasks involving noisy automatic speech recognition, substituting the prompt with "What spoken text can be heard?". The ground truth labels are then the objects within the transcriptions.

Metrics: We propose two metrics for evaluating generative audio captioning tasks. First, similar to the concept of Caption Hallucination Assessment with Image Relevance (CHAIR) [38], we propose the metric named ECHO, which is Evaluation of Caption Hallucination in audiO, to evaluate object hallucination in audio captioning tasks. Given the ground truth objects in the audio, ECHO calculates the proportion of objects appearing in the caption but not being present in the audio. We adopt two variants, ECHO_I and ECHO_S, to evaluate hallucination at the object instance and sentence levels, respectively. They can be expressed as:

$$ECHO_I = \frac{|\{\text{hallucinated objects}\}|}{|\{\text{mentioned objects in the caption}\}|}, \quad (1)$$

$$ECHO_S = \frac{|\{\text{captions with hallucinated objects}\}|}{|\{\text{all captions}\}|}. \quad (2)$$

Second, Cover (Cov) measures the extent to which response cover the audio captions. It can be expressed as:

$$Cover = \frac{|\{objects \text{ in the caption that actually exist in audio}\}|}{|\{ground \text{ truth objects}\}|}.$$
(3)

Table 1: Results of discriminative tasks for audio captioning: Acc (Accuracy), P (Precision), R (Recall), F1 (F1 scores), and Std (standard deviation of F1 scores across five prompts). Greedy and Sample refer to the decoding strategies. (Unit: %)

POPE	Model	Sample						Greedy					
		Acc	P	R	F1	Yes	Std	Acc	P	R	F1	Yes	Std
	Qwen-Audio-Chat-7B	65.3	79.2	32.8	46.1	79.3	12.9	65.0	97.0	30.5	46.4	84.3	23.0
	LTU-AS-7B	50.1	49.2	46.5	47.8	52.8	5.2	51.7	52.0	43.9	47.7	57.8	16.8
Random	SALMONN-7B	56.3	90.0	14.1	24.4	92.2	29.5	56.0	89.9	13.4	23.4	92.5	29.8
	SALMONN-13B	63.7	95.7	28.6	44.1	85.1	25.7	67.6	96.6	36.6	53.1	81.1	15.3
	Specialized-LLaMA	65.5	60.9	82.7	70.1	32.0	3.1	67.5	61.9	88.6	72.9	28.4	3.6
	Specialized-ChatGPT	77.1	69.3	96.6	80.7	30.3	0.7	79.8	71.9	97.6	82.8	32.1	0.7
	Owen-Audio-Chat-7B	59.3	70.4	20.8	32.2	85.2	13.4	58.1	96.8	16.8	28.4	90.7	25.3
	LTU-AS-7B	47.1	45.6	40.5	42.9	55.6	5.2	47.3	46.4	35.0	39.9	62.2	14.4
Popular	SALMONN-7B	56.8	90.4	15.2	26.0	91.6	27.3	65.0	95.4	31.6	47.4	83.4	27.6
Adversarial	SALMONN-13B	65.2	96.0	31.7	47.6	83.5	23.4	65.2	96.1	31.6	47.6	83.6	23.6
	Specialized-LLaMA	58.1	55.4	68.0	61.1	38.7	5.0	60.2	57.0	74.0	64.4	35.0	2.0
	Specialized-ChatGPT	66.2	64.7	70.6	67.5	45.4	0.8	66.5	65.1	70.9	67.8	45.5	0.6
	Qwen-Chat-7B	57.6	70.7	16.6	26.8	88.4	12.4	56.0	89.5	11.8	20.8	93.5	19.7
	LTU-AS-7B	49.0	47.4	44.2	45.8	53.8	4.7	50.1	49.6	40.6	44.6	59.5	15.8
	SALMONN-7B	55.2	87.5	11.2	19.9	93.6	22.4	60.3	93.3	21.3	34.7	88.7	22.7
	SALMONN-13B	60.4	94.3	21.3	34.7	88.8	18.3	60.4	94.3	21.3	34.8	88.8	18.5
	Specialized-LLaMA	59.8	56.6	71.6	63.2	37.3	5.2	62.7	58.6	79.3	67.4	32.9	1.6
	Specialized-ChatGPT	70.4	66.6	79.2	72.3	41.1	0.4	70.7	67.1	79.6	72.8	41.2	0.4

Table 2: Results of generative tasks. E stands for ECHO, with subscripts I and S indicating instance and sentence level, respectively. The subscript g refers to the use of GPT-4 as the evaluator. Std denotes the standard deviation of ECHO across five prompts. SPICE is referenced from [34]. Due to its reliance on fixed template inputs, Qwen-Audio cannot accommodate varying prompts. Specialized in audio captioning refers to specialized audio caption model [35], while in ASR, it refers to Whisper [36] large v3. (Unit: %)

			Audio Captioning								Noisy Speech Recognition			
Strategy	Model	$E_I \downarrow$	$E_S \downarrow$	Cov ↑	StdI	Std_S	$E_{I,g} \downarrow$	$E_{S,g} \downarrow$	Cov _g ↑	SPICE ↑	$E_I \downarrow$	$E_S \downarrow$	Cov ↑	WER ↓
	Qwen-Audio-Chat-7B	39.5	66.3	13.1	0.9	1.4	25.7	44.0	14.6	22.2	25.3	76.6	56.6	44.0
Sample	Qwen-Audio-7B	38.0	58.6	11.5	-	-	25.8	45.0	3.4	11.1	17.3	77.6	36.3	30.3
	LTU-AS-7B	85.2	90.0	2.1	8.7	9.5	84.1	93.9	0.8	5.3	47.6	53.2	56.1	97.0
	SALMONN-7B	33.5	57.1	15.2	15.2	19.2	20.7	36.7	18.8	23.4	30.9	62.6	44.6	73.3
	SALMONN-13B	43.6	69.1	12.7	15.6	21.4	22.7	40.2	19.2	21.5	32.9	61.0	46.6	67.0
	Specialized	33.5	56.5	13.2	-	-	26.2	44.6	15.4	13.2	16.8	33.2	74.9	37.0
Greedy	Owen-Audio-Chat-7B	27.8	50.3	15.1	2.1	5.7	20.4	35.0	16.7	22.2	18.5	42.4	80.0	31.0
	Owen-Audio-7B	25.4	38.5	12.5	-	-	20.6	33.9	17.1	13.7	14.2	31.1	80.2	41.0
	LTU-AS-7B	91.9	98.9	2.4	4.1	3.1	81.5	91.0	1.7	6.0	14.8	27.3	79.3	32.0
	SALMONN-7B	33.1	54.9	13.9	17.0	24.7	21.9	37.6	18.8	21.87	61.0	68.6	34.5	73.0
	SALMONN-13B	43.8	69.3	12.4	15.7	23.0	22.4	40.0	19.4	21.1	20.9	34.4	68.7	41.0
	Specialized	27.4	47.7	15.0	-	-	24.5	40.5	18.2	16.5	10.4	24.3	81.1	23.0

In order to mitigate potential errors introduced by automated measurement methods, we introduce GPT-4 as a reference baseline. We feed captions generated by LALMs, along with ground truth captions and labels, into GPT-4, requesting it to analyze which objects are considered hallucinations and which are not. Consequently, we can also use the decomposed results of GPT-4 and calculate the ECHO and Cover scores. We name it as ECHO_{I, g}, ECHO_{S, g} and Cover_g. For noisy automatic speech recognition task, we also report the word error rate (WER) as a reference and select Whisper [36] large v3 as a baseline model for comparison.

2.4. Evaluation Settings

We select five publicly available LALMs: Qwen-Audio [3], Qwen-Audio-Chat [3], LTU-AS-7B [1], SALMONNN-7B [4] and SALMONNN-13B [4]. We explore greedy and sample decoding strategies, setting the temperature to 1.0, top p to 0.9, and top k to 50 for sample-based decoding. Sample decoding is performed three times to calculate the average results. In addition, as an intuitive baseline, we cascade an existing Whisper-based audio captioning model [35] with ChatGPT (gpt-3.5-turbo-0125) or LLaMA-7b-chat [39] to serve as a reference, which is named as Specialized-ChatGPT and Specialized-

LLaMA, respectively. For cascade pipeline, we combine captions obtained from the audio captioning model and transcriptions acquired through Whisper with the original discriminative questions as input to LLMs to obtain responses. Our codes are available at: github.com/kuan2jiu99/audio-hallucination.

3. Evaluation Results

3.1. Results on Discriminative Tasks

Illustrated in Table 1, we observe that the recall scores of all LALMs are significantly lower than their precision scores. This suggests that LALMs tend to provide affirmative answers when addressing the issue of hallucination, indicating that LALMs are easily misled by non-existent objects. Given that we maintained a 1:1 ratio between ground truth and non-existent objects, the results from the yes rate reveal that most LALMs are predisposed to giving a "Yes" response.

Second, due to the differences in sampling strategies, the difficulty of answering questions varies for the models, resulting in most LALMs experiencing a decrease in F1 scores according to the Random, Popular, and Adversarial settings. From this, it is evident that LALMs are more prone to hallucinating about objects that appear frequently or concurrently. Additionally, in

Table 3: The F1 scores for discriminative tasks, under various prefix prompts and a random sampling strategy, are reported as differences relative to the baseline. This baseline is determined by the F1 scores from tests conducted without the addition of prefix prompts. P1: "Listen.", P2: "Listen closely to the audio before answering the following question.", P3: "Carefully consider the question before providing your answer.", P4: "Listen closely to the audio and carefully consider the question before providing your answer.", P5: "Focus and answer the following question.", P6: "Focus on the given audio and answer the following question.", P7: "Focus on the question and provide the answer." (Unit: %)

Model	P1	P2	Р3	P4	P5	P6	P7	P8
Qwen-Audio-Chat-7B	-11.2	+5.7	+12.8	+9.0	+6.1	+12.1	0.0	+6.8
LTU-AS-7B	-9.4	-8.2	-5.8	-16.5	-6.8	-12.3	-7.1	-8.6
SALMONN-7B	+5.3	+14.7	+28.4	+30.4	+14.7	+32.0	+22.5	+32.2
SALMONN-13B	+3.7	+8.6	+27.0	+25.1	+7.0	+17.4	+12.9	+5.7

Table 1, we report the standard deviation of the F1 scores for different prompts, highlighting the models' sensitivity to various prompt designs. We also discover that the prompt "Have you noticed the sound of [object]?" consistently yields notably high F1 score performance across all models. Third, F1 score of the cascade pipeline significantly surpasses that of all LALMs, indicating a substantial gap between current LALMs and cascade pipelines that needs to be bridged.

3.2. Results on Generative Tasks

Illustrated in Table 2, the performance of LALMs on the ECHO and Cover metrics is comparable to that of Whisper-based caption model [35]. The results obtained using GPT-4 as evaluator exhibit consistent trends. This indicates that LALMs are capable of generating high-quality audio captions, demonstrating sufficient ability to understand audio information. Furthermore, the extent of object hallucination is similar among them except for LTU-AS, suggesting that LALMs can match the performance of Whisper-based audio captioning models in both understanding audio content and the level of object hallucination. Compared to the significant performance gap on discriminative tasks between LALMs and cascade pipelines, it is evident that even though LALMs exhibit performance that rivals specialized caption models in audio captioning tasks, they falter significantly when faced with discriminative questions. This discrepancy highlights a specific weakness of LALMs in handling tasks that require precise discrimination based on audio content, despite their competency in generating descriptive captions from audio inputs. Since LALMs are capable of understanding information within audio, it suggests that the issue may not lie with their ability to process audio content. Instead, it is likely that LALMs struggle to fully comprehend the nature of discriminative questions, making it challenging for them to extract the required information from the given audio. Besides, compared to discriminative tasks, the greedy decoding strategy effectively reduces hallucination phenomena in generative tasks.

3.3. Results on Prompt Engineering

Building on the observation that LALMs are indeed capable of understanding audio information, suggesting their struggle is not with processing audio content but with fully comprehending the nature of discriminative questions, we add appropriate prefix prompts before discriminative questions. These prefix prompts are denoted as P1 to P8, as shown in Table 3, where F1 scores are reported as differences relative to the baseline F1 score. The baseline score is the result without the addition of prefix prompts. By doing so, we anticipate that LALMs thereby focus on seeking the specific information required from the au-

dio to address the question. For P1 to P4, by emphasizing the importance of listening to the given audio, contemplating the question, or both, we experimented with the impact on LALMs' performance. In Table 3, we discover that instructing LALMs to carefully consider the question and to both attentively listen to the audio and carefully consider the question yielded the most significant improvement in performance except for LTU-AS. Merely emphasizing listening to the given audio do not achieve as substantial an improvement as instructing the model to carefully consider the question. On the other hand, for P5 to P8, by emphasizing focus on the given audio, the question, or both, we examine their effects on LALMs' performance. Our findings indicate that instructing LALMs to concentrate on specific information or both can enhance performance except for LTU-AS. This suggests the importance of clearly directing LALMs towards the focus of their attention. Additionally, we experimented with substituting "Focus" with synonyms like "Pay attention" and "Concentrate", yielding similar results. Conclusively, designing an appropriate prompt significantly impacts the performance on discriminative questions.

4. Conclusion and Future Work

Despite recent advancements, LALMs present reliability concerns, particularly with object hallucination. We deploy task-oriented methodologies to gauge this issue, revealing that despite LALMs' capability to perform audio captioning comparably to specialized models, they struggle significantly with discriminative tasks and exhibit severe object hallucination. Compared to cascade pipelines, LALMs still have considerable ground to cover in addressing object hallucination challenges. We propose improvement methods, whereby instructing LALMs to focus on specific information, listen to the audio, or carefully consider the question before responding can enhance model performance. Future work will refine prompts for LALMs to improve audio information extraction and response accuracy, and develop strategies to lessen hallucination in both pre-training and inference stages.

5. Limitations

In generative tasks, combining with automatic segmentation may lead to mismatches between extracted objects and human annotations, causing discrepancies. Challenges also arise in noun extraction when LALMs inaccurately generate captions. To mitigate these issues, we also utilize GPT-4 as our evaluator. In discriminative tasks, LALMs occasionally ignore instructions, complicating automated evaluation by not providing binary answers.

6. Acknowledgement

We would like to thank En-Pei Hu and Ke-Han Lu for providing valuable feedback on the draft of this paper. Additionally, we thank the National Center for High-performance Computing (NCHC) of National Applied Research Laboratories (NARLabs) in Taiwan for providing computational and storage resources.

7. References

- Y. Gong, A. H. Liu, H. Luo, L. Karlinsky, and J. Glass, "Joint audio and speech understanding," in 2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU), 2023.
- [2] Y. Gong, H. Luo, A. H. Liu, L. Karlinsky, and J. Glass, "Listen, think, and understand," arXiv:2305.10790, 2023.
- [3] Y. Chu, J. Xu, X. Zhou, Q. Yang, S. Zhang, Z. Yan, C. Zhou, and J. Zhou, "Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models," arXiv:2311.07919, 2023.
- [4] C. Tang, W. Yu, G. Sun, X. Chen, T. Tan, W. Li, L. Lu, Z. Ma, and C. Zhang, "Salmonn: Towards generic hearing abilities for large language models," arXiv:2310.13289, 2023.
- [5] Z. Kong, A. Goel, R. Badlani, W. Ping, R. Valle, and B. Catanzaro, "Audio flamingo: A novel audio language model with few-shot learning and dialogue abilities," arXiv:2402.01831, 2024.
- [6] M. Shukor, C. Dancette, A. Rame, and M. Cord, "Unified model for image, video, audio and language tasks," arXiv:2307.16184, 2023
- [7] Y. Shu, S. Dong, G. Chen, W. Huang, R. Zhang, D. Shi, Q. Xiang, and Y. Shi, "Llasm: Large language and speech model," arXiv:2308.15930, 2023.
- [8] H. Zhang, X. Li, and L. Bing, "Video-llama: An instruction-tuned audio-visual language model for video understanding," in Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, 2023, pp. 543– 553.
- [9] C. Lyu, M. Wu, L. Wang, X. Huang, B. Liu, Z. Du, S. Shi, and Z. Tu, "Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration," arXiv:2306.09093, 2023.
- [10] Q. Chen et al., "Lauragpt: Listen, attend, understand, and regenerate audio with gpt," arXiv:2310.04673, 2023.
- [11] S. Liu, A. S. Hussain, C. Sun, and Y. Shan, "Music understanding llama: Advancing text-to-music generation with question answering and captioning," arXiv:2308.11276, 2023.
- [12] J. Wu et al., "On decoder-only architecture for speech-to-text and large language model integration," in 2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU). IEEE, 2023, pp. 1–8.
- [13] S. Deshmukh, B. Elizalde, R. Singh, and H. Wang, "Pengi: An audio language model for audio tasks," *Advances in Neural Infor*mation Processing Systems, vol. 36, 2024.
- [14] M. Wang et al., "Slm: Bridge the thin gap between speech and text foundation models," in 2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU). IEEE, 2023, pp. 1–8.
- [15] J. Pan, J. Wu, Y. Gaur, S. Sivasankaran, Z. Chen, S. Liu, and J. Li, "Cosmic: Data efficient instruction-tuning for speech in-context learning," arXiv:2311.02248, 2023.
- [16] S. Wu, H. Fei, L. Qu, W. Ji, and T.-S. Chua, "Next-gpt: Any-to-any multimodal llm," arXiv:2309.05519, 2023.
- [17] J. Zhan et al., "Anygpt: Unified multimodal llm with discrete sequence modeling," arXiv:2402.12226, 2024.
- [18] C.-y. Huang, K.-H. Lu, S.-H. Wang, C.-Y. Hsiao, C.-Y. Kuan, H. Wu, S. Arora, K.-W. Chang, J. Shi, Y. Peng *et al.*, "Dynamic-superb: Towards a dynamic, collaborative, and comprehensive instruction-tuning benchmark for speech," *arXiv*:2309.09510, 2023

- [19] Q. Yang et al., "Air-bench: Benchmarking large audio-language models via generative comprehension," arXiv:2402.07729, 2024.
- [20] R. OpenAI, "Gpt-4 technical report," arXiv, pp. 2303–08774, 2023.
- [21] Y. Zhang et al., "Siren's song in the ai ocean: a survey on hallucination in large language models," arXiv:2309.01219, 2023.
- [22] L. Huang et al., "A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions," arXiv:2311.05232, 2023.
- [23] V. Rawte, A. Sheth, and A. Das, "A survey of hallucination in large foundation models," arXiv:2309.05922, 2023.
- [24] Y. Li, Y. Du, K. Zhou, J. Wang, W. X. Zhao, and J.-R. Wen, "Evaluating object hallucination in large vision-language models," in Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, 2023, pp. 292–305.
- [25] J. Wang et al., "Evaluation and analysis of hallucination in large vision-language models," arXiv:2308.15126, 2023.
- [26] Y. Zhou, C. Cui, J. Yoon, L. Zhang, Z. Deng, C. Finn, M. Bansal, and H. Yao, "Analyzing and mitigating object hallucination in large vision-language models," in *NeurIPS 2023 Workshop on In*struction Tuning and Instruction Following, 2023.
- [27] J. Wang, Y. Wang, G. Xu, J. Zhang, Y. Gu, H. Jia, M. Yan, J. Zhang, and J. Sang, "An Ilm-free multi-dimensional benchmark for mllms hallucination evaluation," arXiv:2311.07397, 2023.
- [28] H. Lovenia, W. Dai, S. Cahyawijaya, Z. Ji, and P. Fung, "Negative object presence evaluation (nope) to measure object hallucination in vision-language models," arXiv:2310.05338, 2023.
- [29] W. Dai, Z. Liu, Z. Ji, D. Su, and P. Fung, "Plausible may not be faithful: Probing object hallucination in vision-language pretraining," in *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, 2023, pp. 2136–2148.
- [30] B. Zhai, S. Yang, C. Xu, S. Shen, K. Keutzer, and M. Li, "Halle-switch: Controlling object hallucination in large vision language models," arXiv e-prints, pp. arXiv-2310, 2023.
- [31] C. D. Kim, B. Kim, H. Lee, and G. Kim, "Audiocaps: Generating captions for audios in the wild," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019, pp. 119–132.
- [32] J. F. Gemmeke, D. P. W. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, "Audio set: An ontology and human-labeled dataset for audio events," in 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2017, pp. 776–780.
- [33] S. Watanabe et al., "Chime-6 challenge: Tackling multispeaker speech recognition for unsegmented recordings," in 6th International Workshop on Speech Processing in Everyday Environments (CHiME 2020). ISCA, 2020.
- [34] P. Anderson, B. Fernando, M. Johnson, and S. Gould, "Spice: Semantic propositional image caption evaluation," in *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part V 14.* Springer, 2016, pp. 382–398.
- [35] M. Kadlčík, A. Hájek, J. Kieslich, and R. Winiecki, "A whisper transformer for audio captioning trained with synthetic captions and transfer learning," 2023.
- [36] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *International Conference on Machine Learning*. PMLR, 2023, pp. 28 492–28 518.
- [37] M. Honnibal and I. Montani, "spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing," 2017, to appear.

- [38] A. Rohrbach, L. A. Hendricks, K. Burns, T. Darrell, and K. Saenko, "Object hallucination in image captioning," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 4035–4045.
- [39] H. Touvron *et al.*, "Llama: Open and efficient foundation language models," *arXiv:2302.13971*, 2023.