

2024年8月12日 19:32

Chen Wang^{1,2}, Mingpu Liao¹, Zhongqiang Huang¹, Jinliang Lu^{1,2}, Junhong Wu^{1,2}
Yuchen Liu¹, Chengqiang Zong^{1,2}, Jiajun Zhang^{1,2,4}
¹ Institute of Automation, Chinese Academy of Sciences
² Machine Intelligence Technology Lab, Alibaba Inc.
³ School of Artificial Intelligence, University of Chinese Academy of Sciences
⁴ Wuhan AI Research

wangchen2020@ia.ac.cn (<http://jiazhang.cqzong@nlpr.ia.ac.cn>)
minpeng_lsp_zhijiang9@alibaba-inc.com

论文分析：语音文本模态对齐是gpt-4o关键的一步，本文在输入上做到了对齐，但是副语言信息的缺失意味着无法合成逼真人声，只完成了一部分

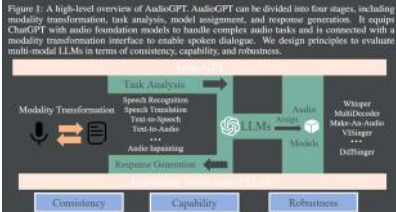
需要考虑的还有 流式实现的问题以及如何获取包含副语言信息特征表征的问题，但论文中迈出了SLLM获取智能这一步，模态对齐和关注续写的训练方法可以借鉴

LLM已经发展到一个很高的水平 -> 语音和文本的天然对齐 -> 如何利用少量的语音数据激发LLM的语言能力 -> 将LLM能力延伸到语音模式

已有的方法:

1. 级联的系统: 比如 audioGPT: speechin/speechout + GPT + audio model

组装起来的用于音频分析生成的系统：能够在多轮对话中处理包含语音、音乐、声音、说话人等理解和生成的AI任务，使人类能够以前所未有的轻松方式创建丰富多样的音频内容。



2. 端到端的语音语言大模型：比如 LLASLM，核心思想是stage1 利用ASR任务搭建简单的asr相关指令（帮我把这句话转成文字）构建模型适应pretrain model; stage2 构建跨模态的语音语言微调指令集进行微调，重点是指令集构建与发布，模型结构与训练过程与本文基本一致。作者认为该文章缺点是：依赖指令集且基于asr任务导致语音特征表现出明显偏差鲁棒性较差

★ mask: <https://huggingface.co/datasets/LinkSoul/LLaSM-Audio-Instructions> -> 公开指令集收集; 需要了解的是指令集构建方法, 以及训练数据的构造



Figure 4: The sample sequence format for the pre-training. We follow the manner of LlaMA-2 and B_INST = [INST] E_INST = [INST], R_SYS = <<SYS>>, I_U = E_SYS + <<I>>, S = <<S>>. The SYSTEM = <<You are a helpful language and speech assistant. You are able to understand the speech content that the user provides, and assist the user with a variety of tasks using natural language.>> and the TEXT_LABEL is the text label of the ASR data associated with the audio input as defined by the definition of the ASR dataset. AUDIO_TOKEN = <<[audio token]>>, AUDIO_PATCH_TOKEN = <<[audio patch token]>>. The *audio_tokens* contains the *audio_tokens*, and the *i_audio_p* is which *i_audio_p* is a single instruction, and is randomly put before or after the *audio_tokens*. While training the BOS tokens and the EOS tokens will be added to each sample at the beginning and the end of the sequence, only the green tokens are used to compute the loss.



Figure 5: The sample sequence format for the cross-modal question answering task. We follow the manner of Liang et al. and BERT+INSTT, i.e., INSTT = [INSTT, B_RVYS] < CRYSTAL > [IAC, E_SVS] + n (< CRYSTAL > -> IAC). The SVS ITEM = "You are a helpful language and speech assistant. You are able to understand the speech content that the user poses questions, and assist the user with a variety of tasks using natural language." and the TEXT_ITEM = "The following text contains information about the user's question. Please provide a detailed answer to the question based on the information provided in the text." The AUDIO_START_TOKEN = "<audio_start>", AUDIO_END_TOKEN = "<audio_end>", AUDIO_PATCH_TOKEN = "<audio_patch>". The control tokens, i.e., the next tokens, which will be replaced by the audio embeddings during training. Each round of question and answer generation is concatenated together with the EOS token. When generating the BGS, the word "BGS:" will be added at the end of the sequence, and the EOS token will be added at the end of the sequence, only the green labels are used to constrain the loss.

Table 1: LLaSM-Audio-Instructions Data

LLaSM-Audio-Instructions				
Source	Conversations	Samples	English Samples	Chinese Samples
WizardLM	80k	160k	159k	<1k
ShareGPT	21k	155k	140k	15k
GPT-4-LLM	96k	192k	128k	64k
Total	199k	508k	428k	80k

本文中，提出了 BLSP 方法，利用现有的 ASR 训练数据，通过行为对引导语音语言预训练。加入冻结的语音 encoder 和 LLM 之间的轻量级模态适配器来实现这一点，确保 LLM 表现出相同忽略模态（语音 or 文本）的生成行为。文中重点关注读写能力，与原有 LLM 的能力相似。文中证明了对引导语音训练的算法可以将 LLM 的能力扩展到语音，并实现与现有系统相比具有竞争力的性能，在零样本跨语言场景也可以做到同样的能力，从而实现语音识别、语音翻译、口语理解和语音对话。

主要贡献:

1. 提出BLSP方法即通过行为对齐实现语言-语音预训练模型的训练, 引入一种新的LLM的多模态对齐方案; 且改进了基于ASR的pretrain改为关注续写任务的pretrain, 更贴合LLM本身的能力
2. 轻量级模块适配器, 训练好的encoder和LLM中间加入适配器做到少量参数训练, 多模态适配; 这样通过asr训练数据就可以做到模态适配, 减少了对语音指令数据的依赖
3. 丰富的实验验证和结果分析

通过行为对齐的方法进行语言语音模型的预训练:

发现分析问题: 对几个不同的任务 续写任务, 情绪分析, 语音识别, 语音翻译任务的 不同模态 (语音 or 文本) 输入的特征表征进行分析, 直觉上, 语音对齐LLMToken上映射的好的话, 那么相同任务的不同模态的特征表征应该表现出靠近, 而不同任务的相同或者不同模态都应该分离

基于ASR pretrain model的模态特征表现出了，语音表征跟语音表征在一起玩，且完全不跟指令去跟相应的text表征合作的迹象，这不是好的模态融合结果，也就是行为跟表征没有对齐，各干各的/

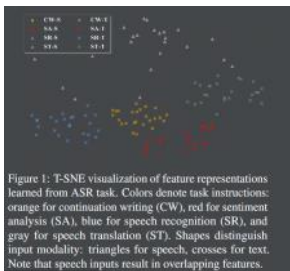


Figure 1: T-SNE visualization of feature representations learned from ASR task. Colors denote task instructions: orange for continuation writing (CW), red for sentiment analysis (SA), blue for speech recognition (SR), and gray for speech translation (ST). Shapes distinguish input modality: triangles for speech, crosses for text. Note that speech inputs result in overlapping features.

	CW-S	SA-S	SR-S	ST-S
CW-S	1.000	0.997	0.997	0.991
SA-S	0.997	1.000	0.997	0.992
SR-S	0.997	0.997	1.000	0.993
ST-S	0.991	0.992	0.993	1.000

Table 7: Average similarity between representations of the same speech inputs under different task instructions learned from ASR task.

CW	SA	SR	ST
0.270	0.106	0.328	0.176

Table 8: Average similarity between representations of paired speech/text inputs under the same task instructions learned from ASR task.

解决: 提出**行为对齐**概念, 让模型不再分开关注语音文本表征, 而是关注指令和输出的对齐, 忽视输入模态的不一样。具体分为两步:
step1:

sists of two steps. In the first step, we use speech transcripts as input and instruct the LLM to generate responses using the following prompt:

```
###[Human]:<instruction><transcript>\n\n\n###[Assistant]:
```

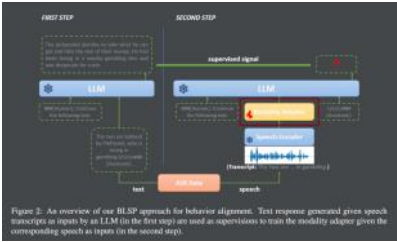
step2:
transcript改成对应的speech signal, step1生成的结果作为ground truth进行有监督训练

sponses produced in the first step as the ground truth for supervised learning using language modeling loss with the following prompt:

```
###[Human]:<instruction><speech>\n\n\n###[Assistant]:<response>
```

此外:在几个任务中,训练时重点关注的续写任务,避免了asr预训练中容易出现的过拟合问题

关于adapter, 是比较常见的模态融合方法, 修改参数量少, 好收敛。冻结encoder和LLM让adapter学习一个模态表征映射就可以



实验:

论文提供了两个版本模型: BLSP 和BLSP-RP 一个只关注续写一个关注续写和识别数据比例9: 1

主要关注的对比模型是CTC+LLM, 也是encoder和LLM的组合

1. 在各个不同任务上的结果:

Method	ASR (WER↓)				ST (BLEU↑)				SLU (ACC↑)			
	LibriSpeech	YELP-LJM	MUSDB18	2	CoVAST	SNIPS	High class	FSC	SLUR	WordSub		
Text+LLM	0.0	5.6	0.0	14.5	19.7	21.8	86.5	52.4	70.9			
Whisper+LLM	3.4	8.9	4.1	10.6	16.8	16.7	83.3	58.9	74.1			
CTC+LLM	8.2	10.8	8.4	20.7	13.3	13.2	79.0	80.4	74.1			
ASR pretraining	—	3.7	—	4.5	0.0	0.0	0.0	0.0	0.0			
BLSP	—	3.8	—	23.1	12.7	12.7	72.8	80.9	76.0			
BLSP+RP	—	6.4	—	8.1	14.9	13.8	78.8	77.3	75.3			

Table 2: Overview of BLSP results on zero-shot speech-to-text tasks. For each ASR test set, we report two WER scores: on the left for the standalone ASR component of a cascaded system, and on the right for instructing the LLM to repeat the words recognized by the ASR component. The BLEU scores for the ST test sets are averaged across multiple modulation directions.

2. 然后验证了 基于blsp方法预训练的模型比基于asr预训练模型在下游任务上更鲁棒:

Method	cs-de	cs-en	cs-fr	cs-it	cs-es	cs-pt	cs-ro	cs-cs
wh pretraining	23.1 / 78.4	25.4 / 78.1	20.9 / 75.6	20.6 / 76.1	23.8 / 78.8	23.3 / 78.7	16.4 / 76.7	13.7 / 72.5
ASR pretraining	23.2 / 78.8	22.8 / 78.7	15.1 / 77.7	22.1 / 78.2	23.6 / 78.7	22.1 / 78.6	16.6 / 77.8	14.8 / 76.5
BLSP	23.3 / 77.7	27.4 / 78.5	11.9 / 76.8	23.2 / 78.8	26.4 / 80.8	24.3 / 80.4	19.2 / 76.8	18.4 / 77.3

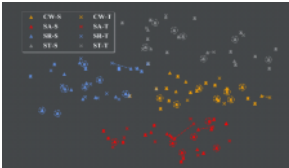
Table 3: ST results (BLEU / COMET) of fine-tuned models on MUST-C.

3. QA 能力验证:

Method	Accept Rate (%)
Text+LLM	94.5
Whisper+LLM	91.3
CTC+LLM	88.6
ASR pretraining	0.0
BLSP	88.5
BLSP+RP	88.3

Table 4: ChatGPT evaluation using acceptable rate.

4. 不同模态输入获取的表征匹配:



在附录 D 中, 我们提供了定量证据, 证明我们的 BLSP 模型可以在不同的指令下为相同的语音输入生成不同的表示, 并且当给出相同的指令时, 对语音和文本输入的表示非常匹配。这些结果表明, BLSP 方法有效地对齐同一空间中的语音和文本输入, 从而将 LLM 的指令跟踪能力扩展到语音输入。

5. 多轮对话能力, 模态可以对齐的话, 多轮对话对LLM来说也不成问题:

