

# VOXDIALOGUE: CAN SPOKEN DIALOGUE SYSTEMS UNDERSTAND INFORMATION BEYOND WORDS?

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

With the rapid advancement of large models, voice assistants are gradually acquiring the ability to engage in open-ended daily conversations with humans. However, current spoken dialogue systems often overlook multi-modal information in audio beyond text, such as speech rate, volume, emphasis, and background sounds. Relying solely on Automatic Speech Recognition (ASR) can lead to the loss of valuable auditory cues, thereby weakening the system’s ability to generate contextually appropriate responses. To address this limitation, we propose **VoxDialogue**, a comprehensive benchmark for evaluating the ability of spoken dialogue systems to understand multi-modal information beyond text. Specifically, we have identified 12 attributes highly correlated with acoustic information beyond words and have meticulously designed corresponding spoken dialogue test sets for each attribute, encompassing a total of 4.5K multi-turn spoken dialogue samples. Finally, we evaluated several existing spoken dialogue models, analyzing their performance on the 12 attribute subsets of VoxDialogue. Experiments have shown that in spoken dialogue scenarios, many acoustic cues cannot be conveyed through textual information and must be directly interpreted from the audio input. In contrast, while direct spoken dialogue systems excel at processing acoustic signals, they still face limitations in handling complex dialogue tasks due to their restricted context understanding capabilities. All data and code will be open source at <https://voxdialogue.github.io/>.

## 1 INTRODUCTION

Voice assistants have rapidly evolved into a focal point of both academic research and industry innovation, aiming to facilitate daily conversations (Li et al., 2017; Lee et al., 2023) and task-oriented dialogues (Budzianowski et al., 2018; Si et al., 2024) with humans. Early iterations relied heavily on automatic speech recognition (ASR) (Yu & Deng, 2016), combined with dialogue understanding and state management, to support basic, predefined tasks. However, these systems (Hoy, 2018) were constrained by their limited scope and inability to handle open-ended interactions. The advent of large language models (LLMs) (Touvron et al., 2023) with enhanced understanding and reasoning capabilities has revolutionized voice assistants, enabling them to engage in more dynamic and unrestricted dialogues with users (OpenAI, 2024b). This marks a significant departure from their earlier, more constrained functionalities, opening up new possibilities for human-computer interaction.

Yet, despite these advancements, current spoken dialogue systems (Zhang et al., 2023; Xie & Wu, 2024; Fang et al., 2024) often overlook the rich multimodal information embedded in audio beyond mere spoken words—such as intonation, volume, rhythm, and background sounds. Relying solely on ASR leads to the omission of valuable auditory cues, diminishing the system’s ability to generate contextually appropriate responses. For example, a system might fail to adjust its language to match a user’s emotional state or regional accent, such as responding with “Yes, madam” to a female voice or adopting british colloquialisms when detecting a British accent.

To address these limitations, recent research has shifted towards developing multimodal audio-language models that enhance system comprehension of audio inputs. Emotion2Vec (Ma et al., 2023), trained on vast emotional speech data, stands as the first high-quality pre-trained model for emotion recognition. Qwen-Audio 1/2 (Chu et al., 2023; 2024) have been trained on extensive datasets encompassing over 30 audio-related tasks, enabling them to understand various au-

Table 1: **Comparison of spoken language and audio comprehension benchmarks in terms of data types and evaluation dimensions.** **SL.** refers to Spoken Language, while **Dlg.** indicates whether the benchmark evaluates on dialogue tasks. **Aud.** represents audio comprehension, and **Mus.** refers to music comprehension. **Speaker Info** includes attributes such as age (Age), gender (Gen), accent (Acc), and language (Lan). **Paralinguistic Info** covers aspects like emotion (Emo), volume (Vol), speech rate (Spd), speech fidelity (Fid), stress (Str), and non-verbal expressions (NVE). <sup>†</sup>Although LeBenchmark includes a small amount of conversational data (29 hours out of 2933 hours), it does not evaluate on the dialogue tasks. <sup>‡</sup>Please note that although AirBench can assess spoken language comprehension, its evaluation of conversational ability (AirBench-Chat) is based on text-based interactions and does not address spoken dialogue capabilities.

Benchmarks	Types		Evaluation Dimensions			
	SL.	Dlg.	Aud.	Mus.	Speaker Info	Paralinguistic Info
SUPERB (Yang et al., 2021)	✓	✗	✗	✗	✗	✓ (Emo)
SLUE (Shon et al., 2022)	✓	✗	✗	✗	✗	✗
LeBenchmark (Evain et al., 2021)	✓	✗ <sup>†</sup>	✗	✗	✗	✓ (Emo)
AF-Dialogue (Kong et al., 2024)	✗	✓	✓	✓	✗	✗
AirBench (Yang et al., 2024)	✗ <sup>‡</sup>	✓	✓	✓	✓ (Age,Gen)	✓ (Emo)
SpokenWOZ (Si et al., 2024)	✓	✓	✗	✗	✗	✗
SD-EVAL (Ao et al., 2024)	✓	✓	✓	✗	✓ (Age,Gen,Acc)	✓ (Emo)
VoxDialogue (ours)	✓	✓	✓	✓	✓ (Age,Gen,Acc,Lan)	✓ (Emo,Vol,Spd,Fid,Str,NVE)

dialogue types—including speech, audio events, and music. Pushing the envelope further, FunAudioLLM (SpeechTeam, 2024) offers full-scene recognition capabilities, detecting non-verbal sounds like laughter and breathing within speech. StyleTalk (Lin et al., 2024b) is the first spoken dialogue system that enables tailoring responses based on contextual emotional information.

As large-scale audio-language models continue to evolve rapidly, the scientific community has increasingly recognized the urgent need for a comprehensive benchmark to effectively evaluate spoken dialogue systems. While some progress has been made, existing benchmarks often exhibit notable shortcomings. For instance, SUPERB (Yang et al., 2021) is the first benchmark specifically designed for spoken language, but it primarily focuses on coarse-grained semantic understanding tasks, overlooking the importance of various acoustic features. Other benchmarks, such as AirBench (Yang et al., 2024) and Audio-Flamingo (Kong et al., 2024), delve deeply into audio understanding, but their dialogue content is limited to the textual modality, making them unsuitable for evaluating spoken dialogue tasks. SpokenWOZ (Si et al., 2024), though valuable for its real human-computer interaction data, is restricted to task-driven dialogues and lacks detailed fine-grained labels. To address more specific attributes of spoken dialogue, SD-EVAL (Ao et al., 2024) shifts the focus to characteristics like gender, age, accent, and emotion, yet its effectiveness is limited by the use of speech utterances that are not derived from dialogue scenarios.

To better benchmark spoken dialogue systems, we analyzed non-textual multimodal acoustic information that may affect dialogue responses, which can be categorized into three main types: speaker information (*age, gender, accent, language*), paralinguistic information (*emotion, volume, speed, fidelity, stress, and various non-verbal expressions*), and background sounds (*audio and music*). In real-world dialogue scenarios, it is crucial to capture not only the semantic content of the speech but also these acoustic cues to generate more appropriate responses. For example, determining the speaker’s age from their vocal tone can help select a suitable form of address. For each of these attributes, we designed the most appropriate spoken dialogue synthesis pipelines. Leveraging the strong inference capabilities of large language models (LLMs) and high-fidelity text-to-speech (TTS) synthesis, we constructed the VoxDialogue benchmark, comprising 12 dialogue scenarios specifically tailored to different acoustic attributes. As shown in Figure 1, to the best of our knowledge, this is the most comprehensive work focusing on acoustic information in spoken dialogue benchmarks. Based on VoxDialogue, we evaluated several existing spoken dialogue systems, comparing the performance of ASR-based dialogue systems and direct dialogue systems across various

108 acoustic-related tasks. The results demonstrate that ASR-based methods are limited in their ability  
109 to understand the diverse acoustic attributes present in spoken dialogues, highlighting the impor-  
110 tance of developing large-scale audio-language models. At the same time, existing direct dialogue  
111 systems (such as Qwen2-Audio) still exhibit limitations in long-context reasoning, indicating the  
112 need for further improvement in their contextual understanding capabilities.

113 All our code and data will be open-sourced. Our main contributions are as follows:  
114

- 115 • We present the first benchmark for evaluating the ability of spoken dialogue systems to  
116 understand acoustic information beyond speech content, VoxDialogue, which integrates  
117 12 acoustic dimensions, including speaker attributes (*age, gender, accent, language*), par-  
118 alinguistic features (*emotion, volume, speed, fidelity, stress, non-verbal expressions*), and  
119 background sounds (*audio, music*).
- 120 • We were the first to develop distinct spoken dialogue data synthesis methods tailored for  
121 different acoustic attributes. This approach enables large-scale synthesis of spoken dia-  
122 logue data, supporting extensive training for spoken dialogue models and endowing them  
123 with more comprehensive acoustic understanding capabilities.
- 124 • We conducted a systematic evaluation of existing spoken dialogue systems, comparing their  
125 performance in terms of understanding acoustic information, supplemented by a qualitative  
126 analysis using a GPT-based metric. Specifically, inspired by the MOS (Mean Opinion  
127 Score) evaluation mechanism, we provided GPT with descriptive criteria corresponding to  
128 different scores, enabling the evaluation model to more accurately assess each response in  
129 terms of both acoustic attributes and content quality.

## 131 2 RELATED WORKS

### 133 2.1 SPOKEN DIALOG SYSTEM

134  
135 With the development of large-scale language models, increasingly powerful audio-language models  
136 have emerged, utilizing extensive training corpora to achieve comprehensive audio understanding.  
137 SpeechGPT (Zhang et al., 2023) integrates discrete speech units into large language models (LLMs),  
138 making it a speech-centric model. Qwen-Audio 1/2 (Chu et al., 2023; 2024) established the first  
139 large-scale, comprehensive audio model for over 30 audio-related tasks, including speech recogni-  
140 tion, speech translation, audio transcription, and audio event detection. Salmonn (Tang et al., 2023)  
141 addresses task complexity in audio models by introducing more intricate story generation tasks.

142 Building on advancements in audio understanding, a series of spoken dialogue models (e.g.,  
143 Qwen-Audio-Chat) have been developed to facilitate more intelligent human-computer interactions.  
144 Audio-Flamingo (Kong et al., 2024) developed a chat model using a text dialogue dataset centered  
145 on audio events, enabling multi-turn, audio-focused text dialogues. StyleTalk (Lin et al., 2024b)  
146 focused on emotional dialogue tasks and introduced the first spoken dialogue model capable of  
147 generating responses with varying emotional tones.

148 However, existing spoken dialogue models (Xie & Wu, 2024; Fang et al., 2024) primarily focus  
149 on understanding speech content and audio information, with a few work specifically dedicated to  
150 comprehending detailed acoustic information within speech. To address this gap, this paper focuses  
151 on 12 acoustic dimensions that could influence dialogue content, including speaker attributes (age,  
152 gender, accent, language), paralinguistic features (emotion, volume, speed, fidelity, stress, non-  
153 verbal expressions), and background sounds (ambient noise, music). Additionally, we developed  
154 a pipeline for synthesizing spoken dialogue data based on these attributes, aiming to enhance the  
155 ability of spoken dialogue models to understand acoustic information beyond the textual content.

### 156 2.2 SPOKEN LANGUAGE BENCHMARK

157  
158 With the rapid development of large-scale audio models (Chu et al., 2024; SpeechTeam, 2024), the  
159 scientific community has increasingly recognized the need for a comprehensive benchmark to eval-  
160 uate spoken dialogue systems. While some progress has been made, many existing benchmarks  
161 still fall short. For instance, SUPERB (Yang et al., 2021) was the first benchmark specifically de-  
signed for spoken language, but it primarily focuses on coarse-grained understanding tasks. Audio-

Table 2: **Dialogue example of VoxDialog**. The **bold black text** highlights the acoustic and background information unique to spoken dialogue, while the **blue text** emphasizes responses that depend on these information.

Attributes	Example
<b>I. Speaker Information</b>	
Gender	A: Excuse me, I’m looking for a reliable smartphone that suits my needs. ( <b>Man Voice</b> ) B: Certainly, <b>Sir!</b> ... excellent performance and a long-lasting battery, perfect for <b>business tasks</b> .
Age	A: What dishes do you recommend? ( <b>Young Voice</b> ) B: Sure! You might enjoy our specialty <b>beef burger</b> , it’s fantastic with <b>fresh juice</b> .
Accent	A: Could you help me with my schedule? ( <b>British Accent</b> ) B: Of course, I’ll assist you in organizing your <b>timetable</b> .
Language	A: 今天天气怎么样? ( <b>Speaking in Chinese</b> ) B: 今天天气晴朗, 气温适中, 适合外出活动。
<b>II. Paralinguistic Information</b>	
Emotion	A: Can you recommend a good movie? ( <b>Sad Tone</b> ) B: I get it. Maybe Soul—it’s emotional and might <b>bring you some comfort</b> .
Volume	A: Can you tell me how to solve this problem? ( <b>Loud Volume</b> ) B: Could you <b>speak more quietly</b> ? People around are studying.
Speed	A: I’d like to get the burger, a soda, and also the pasta, please. ( <b>Speaking Quickly</b> ) B: Sorry, could you say that <b>more slowly</b> ? I didn’t catch everything you said.
Fidelity	A: Can you tell me the time of the meeting? ( <b>Poor Audio Quality</b> ) B: I’m <b>having trouble hearing you</b> . Could you <b>improve the audio quality</b> ?
Stress	A: I really don’t like making sushi. ( <b>Emphasis on “making sushi”</b> ) B: Oh, I see. If you don’t like making sushi, what <b>other types of cooking</b> do you enjoy?
Non-verbal Expressions	A: Could you help me move this box? ( <b>With a sigh before speaking</b> ) B: Are you feeling okay? It seems like <b>you’re really tired</b> . I can take care of it for you.
<b>III. Background Sound</b>	
Audio Events	What was that sound just now? ( <b>Background sound: airplane engine sound, explosion sound</b> ) That was a loud explosion. It sounded like the plane exploded. Hope no one was hurt.
Music	A: Hey, what instrument is this song played on? ( <b>Music: Piano Song, Sad Song</b> ) B: It should be the piano, it sounds so sad.

Flamingo (Kong et al., 2024) explores audio comprehension but is limited to text-based dialogue data, lacking a focus on evaluating spoken dialogue. AirBench (Yang et al., 2024) was the first to propose using GPT-4 (OpenAI, 2024b) to measure spoken dialogue system performance, but its evaluation set remains constrained to a QA interaction. SpokenWOZ (Si et al., 2024) is a large-scale task-oriented dataset that offers real human interaction data, making it valuable for evaluating task-driven dialogue systems. SD-Eval (Ao et al., 2024), which emphasizes acoustic attributes such as gender, age, accent, and emotion, uses raw audio from confessional-style corpora, making it less suitable for conversational scenarios.

However, due to the difficulty of collecting spoken dialogue data in specific scenarios, none of the current benchmarks can effectively evaluate whether spoken dialogue systems can understand various acoustic information beyond text. To address this gap, we developed VoxDialogue, a benchmark created using synthetic data tailored to these acoustic attributes, and evaluated the ability of existing spoken dialogue systems to comprehend such acoustic information.

## 3 VOXDIALOGUE

### 3.1 OVERVIEW

Spoken dialogue systems are typically used in daily dialogues (Lin et al., 2024a). As shown in Table 2, we evaluate the performance of spoken dialogue systems across these three categories in daily

216 dialogue scenarios. Beyond understanding the speech content, spoken dialogue systems must also  
 217 generate the most appropriate responses by considering the speaker’s emotions, gender, and other  
 218 acoustic-related information. Therefore, unlike traditional text-based dialogue benchmarks (Li et al.,  
 219 2017), we systematically analyze the acoustic characteristics that may influence response content  
 220 and have developed a tailored evaluation set specifically for spoken dialogue systems. The evalua-  
 221 tion set for daily dialogue is divided into the following categories: **I. Speaker Information.** (1) *Age*:  
 222 Responses should be tailored to the speaker’s age, adjusting salutations (e.g., Mrs./Miss) or suggest-  
 223 ing content appropriate for their age group. (2) *Gender*: Responses should be gender-specific, mod-  
 224 ifying salutations (e.g., Mr./Mrs.) or offering preferences based on gender. (3) *Accent*: Responses  
 225 should account for the speaker’s accent, selecting vocabulary that aligns with their speech (e.g.,  
 226 British people may be more accustomed to using ‘timetable’ instead of ‘schedule’). (4) *Language*:  
 227 Responses should be adapted to the speaker’s language, choosing the most appropriate language for  
 228 the response. **II. Acoustic Information.** (5) *Emotion*: Responses should detect the speaker’s emo-  
 229 tional state and provide a suitable reply (e.g., suggesting comforting music when sensing distress).  
 230 (6) *Volume*: Responses should consider the speaker’s volume, asking them to lower or raise their  
 231 voice (e.g., requesting quieter speech in quiet environments). (7) *Speed*: Responses should adjust  
 232 to the speaker’s speech rate, asking them to slow down or clarify if speaking too quickly for com-  
 233 prehension. (8) *Fidelity*: Responses should detect poor audio quality and ask the speaker to repeat  
 234 or improve the clarity of their speech for better understanding. (9) *Stress*: Responses should recog-  
 235 nize emphasis on specific words and tailor replies to focus on the stressed content. (10) *Non-verbal*  
 236 *Expressions*: Responses should account for non-verbal cues such as sighs, detecting emotions like  
 237 tiredness or frustration, and offering assistance accordingly. **III. Background Sound.** (11) *Au-*  
 238 *dio Event*: Responses should recognize relevant audio events and adapt accordingly. (12) *Music*:  
 Responses should adjust to the type and mood of the background music.

### 239 3.2 SPOKEN DIALOGUE GENERATION

241 **Stage1: Dialogue Script Synthesis.** Building on the methodology of previous studies (Lin et al.,  
 242 2024a), we employed large language models with advanced reasoning capabilities to synthesize  
 243 spoken conversation scripts tailored to diverse scenarios and acoustic conditions. Specifically, we  
 244 utilized GPT-4o (OpenAI, 2024a) to pre-generate several rounds of historical conversations, fol-  
 245 lowed by the generation of contextually appropriate responses under various controlled acoustic  
 246 conditions. This approach ensures that the synthesized dialogue scripts capture a wide range of  
 247 acoustic features, thereby enhancing their robustness and diversity.

248 **Stage2: Spoken Dialogue Generation.** In line with previous works (Ao et al., 2024; Lin et al.,  
 249 2024a), we utilized high-fidelity TTS (Du et al., 2024) to generate spoken dialogues correspond-  
 250 ing to the dialogue scripts. We carefully tailored the most appropriate speech synthesis method for each  
 251 attribute during the generation process: **(1) Gender, Speed and Emotion.** We use COSYVOICE-  
 252 300M-INSTRUCT<sup>1</sup> to achieve condition speech generation based on gender and emotion by ad-  
 253 justing style instructions. **(2) Stress, Language, and Non-verbal Expressions.** We achieved  
 254 control over these aspects by adjusting the text content in the COSYVOICE-300M-SFT<sup>2</sup>, adding  
 255  $\langle stress \rangle \langle /stress \rangle$ ,  $[laughter]$ , or changing the language of the text. **(3) Volume, Fidelity,**  
 256 **Audio Events, and Music.** We used COSYVOICE-300M-SFT to generate the basic speech, then ap-  
 257 plied post-processing techniques to fine-tune these specific attributes. The details of post-processing  
 258 are shown in Stage 4. **(4) Age.** We randomly selected 1,000 speaker samples of different ages from  
 259 Hechmi et al. (2021) and Tawara et al. (2021) as reference timbres and used COSYVOICE-300M<sup>3</sup>  
 260 for zero-shot TTS synthesis. **(5) Accent.** We used the industrial-grade TTS tool (edge-TTS<sup>4</sup>), which  
 261 offers over 318 timbre references spanning various regions, languages, and genders to achieve pre-  
 262 cise accent generation.

263 **Stage3: Automatic Verification for Spoken Dialogue.** To ensure the quality of the synthesized  
 264 spoken dialogue data, we first employed a pre-trained model to automatically filter out unqualified  
 265 samples, removing those with generation errors and inconsistent timbre. Specifically, we used the  
 266

267 <sup>1</sup><https://huggingface.co/FunAudioLLM/CosyVoice-300M-Instruct>

268 <sup>2</sup><https://huggingface.co/FunAudioLLM/CosyVoice-300M-SFT>

269 <sup>3</sup><https://huggingface.co/model-scope/CosyVoice-300M>

<sup>4</sup><https://github.com/rany2/edge-tts>



Table 3: Detailed statistics of the corresponding subsets of each attribute in VoxDialogue. Gray fonts indicate that samples of this attribute are included in other subsets. IN (India), CA (Canada), ZA (South Africa), GB (United Kingdom), SG (Singapore), US (United States), and AU (Australia). **Turns** represents the total number of turns in each subset, **Dialog.** indicates the number of dialogues in each subset, **Avg** denotes the average number of turns per dialogue in each subset, and **Dur.** refers to the total duration (in hours) of all dialogues in each subset.

Attributes	Categories	Turns	Dialog.	Avg	Dur.
<b>I. Speaker Information</b>					
Gender	Male, Female	2040	340	6.0	3.17
Age	Youth (15-30), Middle-Aged (30-60), Elderly (60+)	3096	447	6.9	6.05
Accent	IN, CA, ZA, GB, SG, US, AU	1440	240	6.0	2.20
Language	Chinese, English	2892	482	6.0	3.51
<b>II. Paralinguistic Information</b>					
Emotion	Neutral, Happy, Sad, Angry, Surprised, Fearful, Disgusted	1980	330	6.0	2.41
Volume	Loud Volume, Low Volume, Normal Volume	1824	304	6.0	2.08
Speed	High Speed, Low Speed, Normal Speed	2184	364	6.0	2.93
Fidelity	Low Fidelity, Normal Fidelity	2196	366	6.0	3.36
Stress	Stress, No Stress	2354	392	6.0	2.51
NVE	Laughter, No Laughter	2046	341	6.0	3.68
<b>III. Background Sound</b>					
Audio	The caption of different audio. (e.g., The wind is blowing and rustling occurs.)	5000	500	10.0	5.25
Music	The aspect list of different music pieces. (e.g., [steeldrum, higher register, amateur recording])	3734	420	8.9	5.42
<b>Overall</b>		<b>30.7K</b>	<b>4.5K</b>	<b>6.8</b>	<b>42.56</b>

we also present a word cloud of VoxDialogue, where it is evident that the dataset primarily consists of daily dialogue, featuring a large number of natural spoken words such as “yeah,” which are representative of daily spoken interactions. This makes it suitable for assessing the performance of spoken dialogue systems in real-world dialogue scenarios. Additionally, the dataset contains numerous acoustically relevant keywords, such as “heard,” “loud,” and “sound,” further supporting the evaluation of acoustic-related aspects of dialogue understanding.

**Distribution of Dialogue Turns and Duration.** All dialogues in our dataset are multi-turn dialogues. In Figure 1 (e), we show the distribution of dialogue turns, with the majority consisting of 6 turns and a maximum of 10 turns. This allows for a comprehensive evaluation of spoken dialogue systems’ ability to understand contexts of varying lengths. In addition, Figures 1 (b) and 1 (c) illustrate the distribution of each turn and the overall dialogue length, respectively, showing that most sentences are approximately 4 seconds long. This implies that the system must understand the context and reason effectively before generating a response.

**Statistics for Subset of Each Attribute.** We present the detailed statistics of each attribute in VoxDialogue in Table 3, covering 35 different categories across 12 attributes. The average number of turns per dialogue exceeds 6, with each attribute containing more than 300 dialogues, ensuring comprehensive reflection of dialogue capabilities.

## 4 BENCHMARK FOR SPOKEN DIALOGUE SYSTEM

### 4.1 TASK DEFINITION

The task of a spoken dialogue system is to generate appropriate responses based on the contextual information from the sequence of human dialogue (e.g., the user’s utterance sequence) and the preceding assistant response sequence, where the total number of dialogue turns is denoted by  $t$ . The

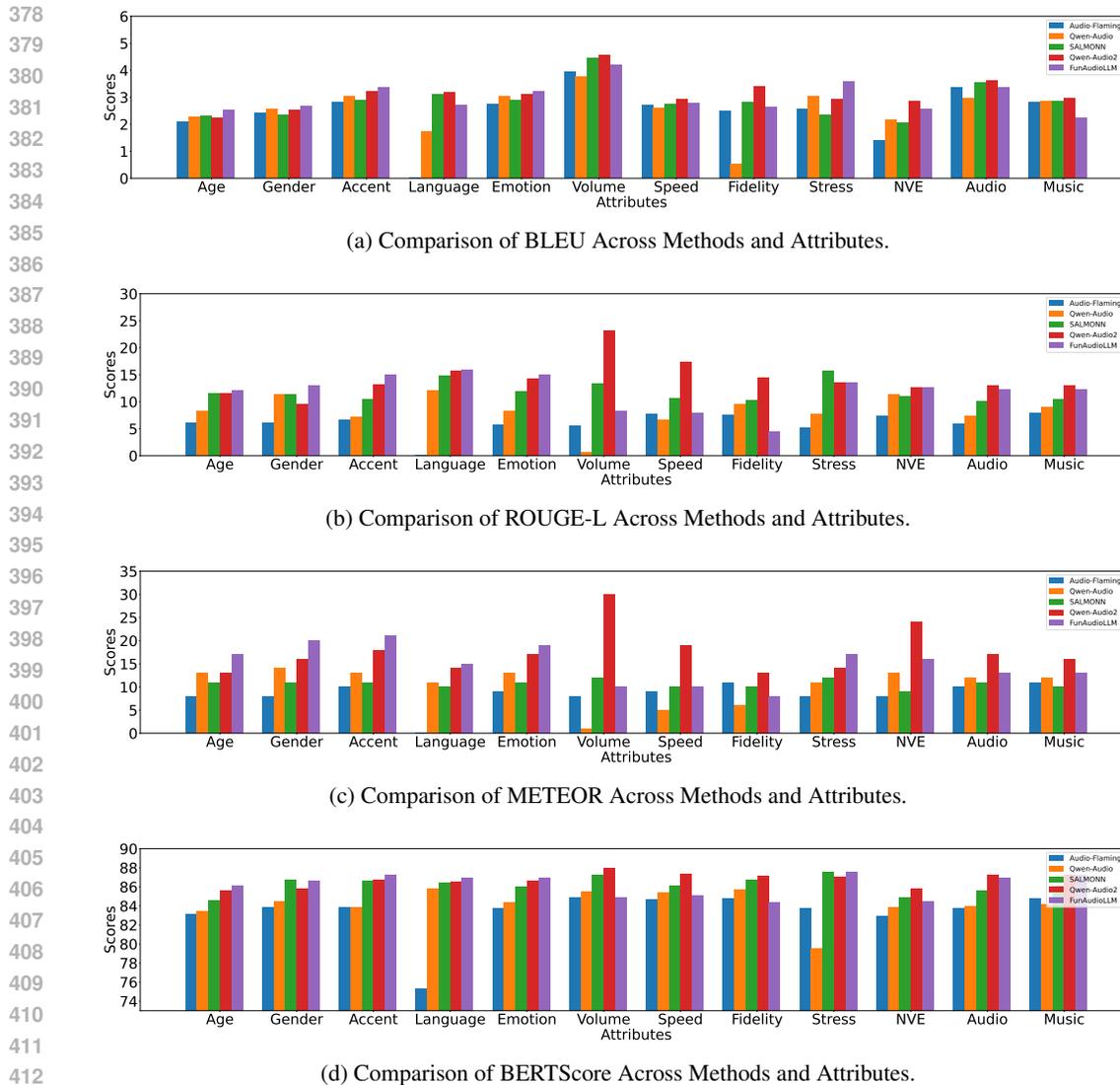


Figure 2: The comparison of spoken dialogue performance across 12 different attribute-specific test sets on the VoxDialogue dataset.

goal of the spoken dialogue system is to generate the most suitable response based on the previous  $t$  utterances and the  $t - 1$  historical replies. In our work, we evaluate the performance of the spoken dialogue system by focusing solely on the final utterance of each dialogue.

## 4.2 EVALUATION METRICS

To assess the model’s performance, we conducted separate tests on a subset of Voxdialogue. Drawing on previous research (Ao et al., 2024), we utilized both quantitative and qualitative metrics for a comprehensive evaluation. The quantitative evaluation focused on two key aspects: content and style. For content evaluation, we employed widely recognized text generation metrics, including vocabulary-level measures such as BLEU (Papineni et al., 2002), ROUGE-L (Lin, 2004), and METEOR (Banerjee & Lavie, 2005), alongside semantic-level metrics like BERTScore (Zhang et al., 2019). For style evaluation, we calculated the weighted F1 score of speech sentiment.

In addition to these quantitative assessments, we conducted a qualitative analysis using GPT-based metric (Yang et al., 2024). The meaning of each score is as follows: **1**: Contextually relevant but lacks attribute information. **2**: Partially relevant to the context but feels unnatural, with no attribute

Table 4: GPT-based Metric Comparison of Different Spoken Dialogue Models on VoxDialogue.

Method	Speaker Info			Paralinguistic Info					Background			
	Age	Gen	Acc	Lan	Emo	Vol	Spd	Fid	Str	NVE	Aud	Mus
<i>ASR-Based Spoken Dialogue System</i>												
FunAudioLLM (SpeechTeam, 2024)	<b>4.32</b>	<b>4.39</b>	<b>3.57</b>	<b>4.61</b>	<b>4.09</b>	1.82	1.92	1.79	3.13	2.87	3.47	3.59
<i>Direct Spoken Dialogue System</i>												
Audio-Flamingo (Kong et al., 2024)	1.00	1.00	1.04	1.72	1.00	1.20	1.14	1.26	1.34	1.06	1.37	1.11
SALMONN (Tang et al., 2023)	1.99	1.64	1.78	3.50	1.84	2.88	2.27	2.29	3.86	2.59	2.15	2.23
Qwen-Audio (Chu et al., 2023)	1.36	1.04	1.28	1.04	1.06	1.48	1.08	1.32	2.49	2.65	1.42	1.18
Qwen2-Audio (Chu et al., 2024)	3.46	4.18	2.71	4.43	3.73	<b>3.06</b>	<b>3.29</b>	<b>2.98</b>	<b>3.93</b>	<b>3.46</b>	<b>3.81</b>	<b>3.98</b>

information. **3**: Partially relevant to the context, with mention of the attribute. **4**: Contextually relevant and natural, mentioning the attribute, but could be improved. **5**: Contextually relevant, smooth, natural, and accurately addresses the attribute. We have included all the evaluated prompt templates in supplementary materials. Please refer to the supplementary materials for more details.

### 4.3 SPOKEN DIALOGUE SYSTEM

In order to build a comprehensive benchmark, we evaluated two main types of spoken dialogue system approaches: (1) **ASR-based dialogue systems** (e.g., FunAudioLLM (Fang et al., 2024)) and (2) **direct spoken dialogue systems**<sup>5</sup> (e.g., Audio-Flamingo (Kong et al., 2024), SALMONN (Tang et al., 2023), Qwen-Audio (Chu et al., 2023), and Qwen2-Audio (Chu et al., 2024)). Figure 2 presents a comparative analysis using four metrics across various attributes on the VoxDialogue dataset. Based on the experimental results, we gained the following key insights:

**ASR-based systems excel in context-sensitive tasks.** In attributes that can be inferred through context understanding, ASR-based systems (such as FunAudioLLM) show significant advantages. ASR systems first transcribe speech into text and then process it, allowing them to more effectively capture and analyze the context of a conversation. For example, in attributes like *Emotion* and *Speaker Information(Age, Gender, Accent, Language)*, FunAudioLLM consistently outperforms direct spoken dialogue systems. The results from BLEU, ROUGE-L, METEOR, and BERTScore metrics indicate that FunAudioLLM achieves higher scores, such as in emotion (3.22 BLEU, 14.93 ROUGE-L, 18.97 METEOR, 86.92 BERTScore). This proves that most current direct spoken dialogue systems lack adequate context understanding capabilities and are far weaker than text-based large language models.

**Advantages of direct spoken dialogue systems in acoustic attribute processing.** Although ASR-based systems can leverage the strong context understanding capabilities of large language models, they struggle with attributes that heavily rely on sound understanding (such as volume, fidelity, speed, and other paralinguistic information). ASR-based methods face challenges when addressing dialogue tasks related to these attributes. In contrast, direct systems like Qwen2-Audio excel in tasks involving these acoustic properties. The results show that Qwen2-Audio outperforms other systems in these categories. For instance, Qwen2-Audio achieved the highest scores for *volume* (4.56 BLEU, 23.13 ROUGE-L, 29.82 METEOR, and 87.98 BERTScore), demonstrating its ability to handle loud and soft speech variations more effectively. Similarly, *fidelity* is another strong point for direct dialogue systems. Qwen2-Audio’s excellent performance in handling varying fidelity levels (3.38 BLEU, 14.36 ROUGE-L, 13.03 METEOR, 87.12 BERTScore) confirms that spoken dialogue tasks, which heavily rely on acoustic information beyond words.

<sup>5</sup>All models used in the evaluation are *-chat* version.

#### 4.4 QUALITATIVE COMPARISON

Inspired by Yang et al. (2024), we also attempted to use GPT-4 (OpenAI, 2024b) for evaluation, focusing on whether the responses exhibit the specific attribute characteristics and whether they provide reasonable replies to the previous context. As shown in Table 4, we present the qualitative testing results of different methods across 12 attributes. Specifically, a score of 3 represents mention of attribute information, 4 represents a reasonable and natural response.

We observed that the conclusions from the qualitative tests largely align with those from the quantitative evaluations. For context-driven attributes (such as speaker information and emotion), ASR-based dialogue models continue to demonstrate the best performance. However, for attributes that are highly dependent on acoustic information (such as speed, fidelity, audio, and music), direct spoken dialogue models like Qwen2-Audio significantly outperform FunAudioLLM, underscoring the importance of developing direct spoken dialogue models.

Additionally, we found that Qwen-Audio often responds with descriptive sentences related to the query, which severely affects its performance. The SALMONN model frequently repeats parts of the query, leading to higher quantitative scores in some attributes (e.g., a BLEUScore of 87.53 for Stress, 0.53 higher than Qwen2-Audio), but its qualitative performance is inferior to Qwen2-Audio (with a GPT-4-based metric score 0.97 lower). This indicates that most current large audio-language models are focused on QA-style interactions, and are not yet well-suited for dialogue-style conversations.

## 5 CONCLUSION

In this work, we introduced **VoxDialogue**, a comprehensive benchmark designed to evaluate spoken dialogue systems’ ability to understand information beyond words. By identifying 12 critical attributes tied to acoustic cues such as speech rate, volume, emphasis, and background sounds, we constructed a challenging test set of 4.5K multi-turn dialogue samples. Our experiments demonstrated that while ASR-based systems excel at context understanding and textual interpretation, they fail to capture important acoustic signals that are essential for contextually appropriate responses. In contrast, direct spoken dialogue systems outperform ASR-based models in processing acoustic properties, but their limited ability to understand complex dialogue contexts remains a significant shortcoming. The findings highlight the importance of acoustic information in enhancing the performance of spoken dialogue systems and reveal the current limitations in both ASR-based and direct spoken dialogue models.

## REPRODUCIBILITY STATEMENT

All of our data, code, and model weights will be open-sourced.

- Section 3 provides detailed instructions on the construction of VoxDialogue, including a comprehensive list of all relevant open-source resources.
- Section 4.1 outlines the detailed task definitions.
- Section 4.2 elaborates on the evaluation metrics and specific details.
- All of our prompt templates are included in the Supplementary Material.

## REFERENCES

- Junyi Ao, Yuancheng Wang, Xiaohai Tian, Dekun Chen, Jun Zhang, Lu Lu, Yuxuan Wang, Haizhou Li, and Zhizheng Wu. Sd-eval: A benchmark dataset for spoken dialogue understanding beyond words. *arXiv preprint arXiv:2406.13340*, 2024.
- Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pp. 65–72, 2005.
- Hervé Bredin. pyannote.audio 2.1 speaker diarization pipeline: principle, benchmark, and recipe. In *Proc. INTERSPEECH 2023*, 2023.

- 540 Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman  
541 Ramadan, and Milica Gasic. Multiwoz-a large-scale multi-domain wizard-of-oz dataset for task-  
542 oriented dialogue modelling. In *Proceedings of the 2018 Conference on Empirical Methods in*  
543 *Natural Language Processing*, pp. 5016–5026, 2018.
- 544 Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and  
545 Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale  
546 audio-language models. *arXiv preprint arXiv:2311.07919*, 2023.
- 547 Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv,  
548 Jinzheng He, Junyang Lin, et al. Qwen2-audio technical report. *arXiv preprint arXiv:2407.10759*,  
549 2024.
- 550 Zhihao Du, Qian Chen, Shiliang Zhang, Kai Hu, Heng Lu, Yexin Yang, Hangrui Hu, Siqi Zheng, Yue  
551 Gu, Ziyang Ma, et al. Cosyvoice: A scalable multilingual zero-shot text-to-speech synthesizer  
552 based on supervised semantic tokens. *arXiv preprint arXiv:2407.05407*, 2024.
- 553 Solène Evain, Ha Nguyen, Hang Le, Marcely Zanon Boito, Salima Mdhaffar, Sina Alisamir, Ziyi  
554 Tong, Natalia Tomashenko, Marco Dinarelli, Titouan Parcollet, et al. Lebenchmark: A repro-  
555 ducible framework for assessing self-supervised representation learning from speech. In *INTER-*  
556 *SPEECH 2021: Conference of the International Speech Communication Association*, 2021.
- 557 Qingkai Fang, Shoutao Guo, Yan Zhou, Zhengrui Ma, Shaolei Zhang, and Yang Feng. Llama-omni:  
558 Seamless speech interaction with large language models. *arXiv preprint arXiv:2409.06666*, 2024.
- 559 Khaled Hechmi, Trung Ngo Trong, Ville Hautamäki, and Tomi Kinnunen. Voxceleb enrichment  
560 for age and gender recognition. In *2021 IEEE Automatic Speech Recognition and Understanding*  
561 *Workshop (ASRU)*, pp. 687–693. IEEE, 2021.
- 562 Matthew B Hoy. Alexa, siri, cortana, and more: an introduction to voice assistants. *Medical refer-*  
563 *ence services quarterly*, 37(1):81–88, 2018.
- 564 Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping, Rafael Valle, and Bryan Catanzaro. Audio  
565 flamingo: A novel audio language model with few-shot learning and dialogue abilities. *arXiv*  
566 *preprint arXiv:2402.01831*, 2024.
- 567 Keon Lee, Kyumin Park, and Daeyoung Kim. Dailytalk: Spoken dialogue dataset for conversational  
568 text-to-speech. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and*  
569 *Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
- 570 Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. Dailydialog: A manually  
571 labelled multi-turn dialogue dataset. In *Proceedings of the Eighth International Joint Conference*  
572 *on Natural Language Processing (Volume 1: Long Papers)*, pp. 986–995, 2017.
- 573 Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization*  
574 *branches out*, pp. 74–81, 2004.
- 575 Guan-Ting Lin, Cheng-Han Chiang, and Hung-yi Lee. Advancing large language models to capture  
576 varied speaking styles and respond properly in spoken conversations. In Lun-Wei Ku, Andre  
577 Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association*  
578 *for Computational Linguistics (Volume 1: Long Papers)*, pp. 6626–6642, Bangkok, Thailand,  
579 August 2024a. Association for Computational Linguistics. URL <https://aclanthology.org/2024.acl-long.358>.
- 580 Guan-Ting Lin, Cheng-Han Chiang, and Hung-yi Lee. Advancing large language models to  
581 capture varied speaking styles and respond properly in spoken conversations. *arXiv preprint*  
582 *arXiv:2402.12786*, 2024b.
- 583 Xubo Liu, Egor Lakomkin, Konstantinos Vougioukas, Pingchuan Ma, Honglie Chen, Ruiming Xie,  
584 Morrie Doulaty, Niko Moritz, Jachym Kolar, Stavros Petridis, et al. Synthvsr: Scaling up visual  
585 speech recognition with synthetic supervision. In *Proceedings of the IEEE/CVF Conference on*  
586 *Computer Vision and Pattern Recognition*, pp. 18806–18815, 2023.

- 594 Ziyang Ma, Zhisheng Zheng, Jiaxin Ye, Jinchao Li, Zhifu Gao, Shiliang Zhang, and Xie Chen.  
595 emotion2vec: Self-supervised pre-training for speech emotion representation. *arXiv preprint*  
596 *arXiv:2312.15185*, 2023.
- 597 OpenAI. Gpt-4o system card. <https://cdn.openai.com/gpt-4o-system-card.pdf>, 2024a.
- 599 OpenAI. Chatgpt can now see, hear, and speak. [https://openai.com/index/chatgpt-can-now-see-](https://openai.com/index/chatgpt-can-now-see-hear-and-speak/)  
600 *hear-and-speak/*, 2024b.
- 602 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic  
603 evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association*  
604 *for Computational Linguistics*, pp. 311–318, 2002.
- 605 Alexis Plaquet and Hervé Bredin. Powerset multi-class cross entropy loss for neural speaker diariza-  
606 tion. In *Proc. INTERSPEECH 2023*, 2023.
- 608 Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever.  
609 Robust speech recognition via large-scale weak supervision. In *International conference on ma-*  
610 *chine learning*, pp. 28492–28518. PMLR, 2023.
- 612 Suwon Shon, Ankita Pasad, Felix Wu, Pablo Brusco, Yoav Artzi, Karen Livescu, and Kyu J Han.  
613 Slue: New benchmark tasks for spoken language understanding evaluation on natural speech. In  
614 *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing*  
615 *(ICASSP)*, pp. 7927–7931. IEEE, 2022.
- 616 Shuzheng Si, Wentao Ma, Haoyu Gao, Yuchuan Wu, Ting-En Lin, Yinpei Dai, Hangyu Li, Rui  
617 Yan, Fei Huang, and Yongbin Li. Spokenwoz: A large-scale speech-text benchmark for spoken  
618 task-oriented dialogue agents. *Advances in Neural Information Processing Systems*, 36, 2024.
- 620 Tongyi SpeechTeam. Funaudiollm: Voice understanding and generation foundation models for  
621 natural interaction between humans and llms. *arXiv preprint arXiv:2407.04051*, 2024.
- 622 Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma,  
623 and Chao Zhang. Salmonn: Towards generic hearing abilities for large language models. *arXiv*  
624 *preprint arXiv:2310.13289*, 2023.
- 626 Naohiro Tawara, Atsunori Ogawa, Yuki Kitagishi, and Hosana Kamiyama. Age-vox-celeb: Multi-  
627 modal corpus for facial and speech estimation. In *ICASSP 2021-2021 IEEE International Con-*  
628 *ference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6963–6967. IEEE, 2021.
- 629 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
630 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
631 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 633 Zhifei Xie and Changqiao Wu. Mini-omni: Language models can hear, talk while thinking in  
634 streaming. *arXiv preprint arXiv:2408.16725*, 2024.
- 635 Qian Yang, Jin Xu, Wenrui Liu, Yunfei Chu, Ziyue Jiang, Xiaohuan Zhou, Yichong Leng, Yuanjun  
636 Lv, Zhou Zhao, Chang Zhou, and Jingren Zhou. AIR-bench: Benchmarking large audio-language  
637 models via generative comprehension. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.),  
638 *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Vol-*  
639 *ume 1: Long Papers)*, pp. 1979–1998, Bangkok, Thailand, August 2024. Association for Compu-  
640 tational Linguistics. URL <https://aclanthology.org/2024.acl-long.109>.
- 642 Shu Wen Yang, Po Han Chi, Yung Sung Chuang, Cheng I Jeff Lai, Kushal Lakhotia, Yist Y Lin,  
643 Andy T Liu, Jiatong Shi, Xuankai Chang, Guan Ting Lin, et al. Superb: Speech processing univer-  
644 sal performance benchmark. In *22nd Annual Conference of the International Speech Communi-*  
645 *cation Association, INTERSPEECH 2021*, pp. 3161–3165. International Speech Communication  
646 Association, 2021.
- 647 Dong Yu and Lin Deng. *Automatic speech recognition*, volume 1. Springer, 2016.

648 Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu.  
649 Speechgpt: Empowering large language models with intrinsic cross-modal conversational abil-  
650 ities. *arXiv preprint arXiv:2305.11000*, 2023.  
651  
652 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluat-  
653 ing text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

Table 5: Detailed Comparison of Spoken Dialogue Systems across Various Metrics

(a) BLEU Scores

Method	Speaker Info				Paralinguistic Info						Background	
	Age	Gen	Acc	Lan	Emo	Vol	Spd	Fid	Str	NVE	Aud	Mus
FunAudioLLM	2.53	2.66	3.34	2.72	3.22	4.20	2.77	2.65	3.58	2.57	3.35	2.25
Audio-Flamingo	2.08	2.40	2.83	0.01	2.74	3.95	2.70	2.50	2.58	1.41	3.38	2.81
Qwen-Audio	2.26	2.56	3.05	1.74	3.01	3.78	2.61	0.54	3.02	2.15	2.97	2.87
SALMONN	2.29	2.35	2.88	3.09	2.88	4.44	2.73	2.82	2.33	2.04	3.55	2.86
Qwen-Audio2	2.22	2.52	3.20	3.18	3.11	4.56	2.92	3.38	2.93	2.85	3.60	2.97

(b) ROUGE-L Scores

Method	Speaker Info				Paralinguistic Info						Background	
	Age	Gen	Acc	Lan	Emo	Vol	Spd	Fid	Str	NVE	Aud	Mus
FunAudioLLM	12.15	12.95	15.07	15.88	14.93	8.28	7.97	4.47	13.49	12.63	12.20	12.20
Audio-Flamingo	6.12	6.15	6.62	0.03	5.78	5.48	7.67	7.57	5.12	7.41	5.91	7.88
Qwen-Audio	8.34	11.44	7.12	12.09	8.24	0.71	6.61	9.58	7.76	11.36	7.29	9.01
SALMONN	11.52	11.43	10.51	14.80	11.81	13.30	10.56	10.22	15.71	11.01	10.05	10.51
Qwen-Audio2	11.51	9.62	13.18	15.66	14.18	23.13	17.34	14.36	13.45	12.67	13.01	12.97

(c) METEOR Scores

Method	Speaker Info				Paralinguistic Info						Background	
	Age	Gen	Acc	Lan	Emo	Vol	Spd	Fid	Str	NVE	Aud	Mus
FunAudioLLM	16.89	20.12	21.03	15.21	19.31	10.19	9.83	8.16	16.95	16.31	13.22	13.04
Audio-Flamingo	8.23	7.79	10.03	0.25	9.17	8.31	8.69	11.04	8.12	7.88	9.93	11.01
Qwen-Audio	12.87	14.16	12.92	11.06	13.12	1.41	5.28	6.11	10.92	13.21	12.22	12.08
SALMONN	11.02	10.81	11.21	10.35	11.13	11.78	10.14	10.17	11.84	9.03	11.18	10.21
Qwen-Audio2	12.96	16.15	18.24	14.37	17.05	30.11	19.08	12.78	14.01	24.41	17.29	15.72

(d) BERTScore

Method	Speaker Info				Paralinguistic Info						Background	
	Age	Gen	Acc	Lan	Emo	Vol	Spd	Fid	Str	NVE	Aud	Mus
FunAudioLLM	86.14	86.65	87.24	86.97	86.87	84.87	85.03	84.36	87.51	84.51	86.98	87.19
Audio-Flamingo	83.12	83.84	83.86	75.28	83.78	84.91	84.71	84.81	83.78	82.89	83.78	84.74
Qwen-Audio	83.41	84.46	83.84	85.79	84.34	85.53	85.34	85.66	79.55	83.85	83.95	84.14
SALMONN	84.64	86.75	86.65	86.44	86.05	87.27	86.06	86.74	87.53	84.92	85.63	86.06
Qwen-Audio2	85.59	85.79	86.67	86.52	86.65	88.02	87.32	87.12	87.08	85.79	87.19	87.19

## A MORE EXPERIMENT RESULTS

### A.1 THE DETAILED PERFORMANCE COMPARISON

For comparison, the detailed performance corresponding to Figure 2 is presented in Table 5.

## B LIMITATION

Our work heavily relies on synthetic datasets. Although prior research (Liu et al., 2023) has shown that synthetic data can be effectively used for training and evaluation, a domain gap persists between

756 synthetic and real-world data. This gap may affect the generalization of models trained on synthetic  
757 data when applied to real-world dialogue scenarios.

758  
759 However, since our focus is on understanding acoustic information, synthetic data proves particu-  
760 larly useful in simulating various acoustic cues found in real conversational settings. Additionally,  
761 the synthetic dataset offers more diverse and controllable dialogue content, making it sufficient for  
762 evaluating whether spoken dialogue systems can understand information beyond text.

763 To properly assess the performance of dialogue systems in real-world scenarios, it is crucial to use  
764 datasets based on authentic conversational environments. We believe that constructing a separate  
765 real-world dialogue evaluation benchmark, independent of our work, would be more effective in  
766 evaluating spoken dialogue systems' performance in real scenarios than using a single dataset to  
767 assess both acoustic information comprehension and real-world dialogue capabilities.

768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809