# OpenS2S: Advancing Fully Open-Source End-to-End Empathetic Large Speech Language Model

Chen Wang[1,2], Tianyu Peng[1,2,3], Wen Yang[1,2], Yinan Bai[1,2],
Guangfu Wang[4], Jun Lin[4], Lanpeng Jia[4],
Lingxiang Wu[1,3], Jinqiao Wang[1,2,3], Chengqing Zong[1,2], Jiajun Zhang[1,2,3] *

[1] Institute of Automation, Chinese Academy of Sciences
[2] School of Artificial Intelligence, University of Chinese Academy of Sciences
[3] Wuhan AI Research  [4] GWM AI Lab
wangchen2020@ia.ac.cn jjzhang@nlpr.ia.ac.cn

## Abstract

Empathetic interaction is a cornerstone of human-machine communication, due to the need for understanding speech enriched with paralinguistic cues and generating emotional and expressive responses. However, the most powerful empathetic LSLMs are increasingly closed off, leaving the crucial details about the architecture, data and development opaque to researchers. Given the critical need for transparent research into the LSLMs and empathetic behavior, we present OpenS2S, a fully open-source, transparent and end-to-end LSLM designed to enable empathetic speech interactions. Based on our empathetic speech-to-text model BLSP-Emo [1], OpenS2S further employs a streaming interleaved decoding architecture to achieve low-latency speech generation. To facilitate end-to-end training, OpenS2S incorporates an automated data construction pipeline that synthesizes diverse, high-quality empathetic speech dialogues at low cost. By leveraging large language models to generate empathetic content and controllable text-to-speech systems to introduce speaker and emotional variation, we construct a scalable training corpus with rich paralinguistic diversity and minimal human supervision. We release the fully open-source OpenS2S model, including the dataset, model weights, pre-training and fine-tuning codes, to empower the broader research community and accelerate innovation in empathetic speech systems.

| | Demo | https://casia-lm.github.io/OpenS2S |
|---|---|---|
| | Code | https://github.com/CASIA-LM/OpenS2S |
| | Model | https://huggingface.co/CASIA-LM/OpenS2S |
| | Data | https://huggingface.co/datasets/CASIA-LM/OpenS2S_Datasets |

## 1 Introduction

Empathy is a fundamental pillar of human interactions, fostering everything from prosocial behavior to deeper connections [2]. Modeling and understanding empathy is a complex task for artificial intelligence, yet its integration is crucial for fostering more natural and effective human-machine

---

* Corresponding author

Technical Report.

Table 1: The degree of openness of Open LSLMs.

| Name | LLaMA-Omni2 | Qwen2-Audio | GLM-4-Voice | Kimi-Audio | OpenS2S |
|---|---|---|---|---|---|
| Training Data | ✗ | ✗ | ✗ | ✗ | ✔ |
| Training Code | ✗ | ✗ | ✗ | ✔ | ✔ |
| Model | ✔ | ✔ | ✔ | ✔ | ✔ |
| Empathetic | ✗ | ✗ | ✔ | ✔ | ✔ |

communication [3]. In the realm of Large Speech Language Models (LSLMs), this challenge is particularly pronounced. Speech inherently conveys a wealth of rich paralinguistic information, including intonation, rhythm, volume variations, and cues related to speaker attributes like age and gender. This intricate paralinguistic content makes speech communication highly sensitive, rendering flat or unnuanced responses from automated systems unacceptable. Consequently, developing empathetic speech systems is vital for creating more natural and human-centered artificial intelligence.

While recent advancements in LSLMs have significantly enhanced audio processing and enabled robust semantic-based instruction following in conversations [4–10], most existing models tend to overlook critical paralinguistic information in speech, thereby fundamentally limiting their native empathetic interaction capabilities. Although some LSLMs [11–13] demonstrate strong empathetic performance, they typically necessitate extensive pre-training on millions of hours of high-quality speech data. This reliance on vast datasets incurs substantial annotation, computation, and training costs, setting up a significant barrier to their broader adoption and development. Furthermore, many of the most advanced models, particularly commercial models like GPT-4o [14] and Gemini are fully proprietary and closed-source. This lack of transparency makes it challenging to analyze their internal mechanisms, reproduce their empathetic behaviors, or build upon their architectures for further scientific research and development. To scientifically study the empathetic behaviors in LSLMs, including potential biases, cultural variations, and their ethical implications, we believe that access to powerful, fully open empathetic LSLMs is critical to the advancement of this field.

To address the aforementioned limitations, we propose OpenS2S, a fully open-source, end-to-end LSLM. OpenS2S not only exhibits competitive foundational speech capabilities but also features an efficient streaming architecture based on interleaved decoding. Crucially, in contrast to existing models that achieve empathetic capabilities through resource-intensive pre-training, OpenS2S attains comparable empathetic interaction performance with significantly lower training data and computational costs. Moreover, the empathetic support provided by OpenS2S transcends mere paralinguistic cues, extending deeply into the semantic content of the dialogue.

Overall, our main contributions are as follows:

1. **Model Construction and Training:** We build an efficient speech-to-speech empathetic model based on an advanced framework and conduct extensive training using high-quality data. This model can provide a more convenient and natural way for humans to interact with artificial intelligence.

2. **Automatic Empathetic Speech Instruction Dataset Construction:** We propose a data augmentation method for empathetic speech dialogue by combining the strengths of large language models (LLMs) and text-to-speech (TTS) models. LLMs are used to generate diverse user queries and empathetic responses, while voice cloning ensures input speaker diversity. InstructTTS further enables controllable emotional expression in speech responses, facilitating the construction of rich, high-quality training data with minimal human annotation.

3. **Fully Open-Source Release:** To foster collaborative research and accelerate innovation in empathetic LSLMs, we release all the resources, including the model weights, all codes for constructing datasets, pre-training, fine-tuning and evaluation, and the synthetic datasets, providing fully transparency and reproducibility for the community.
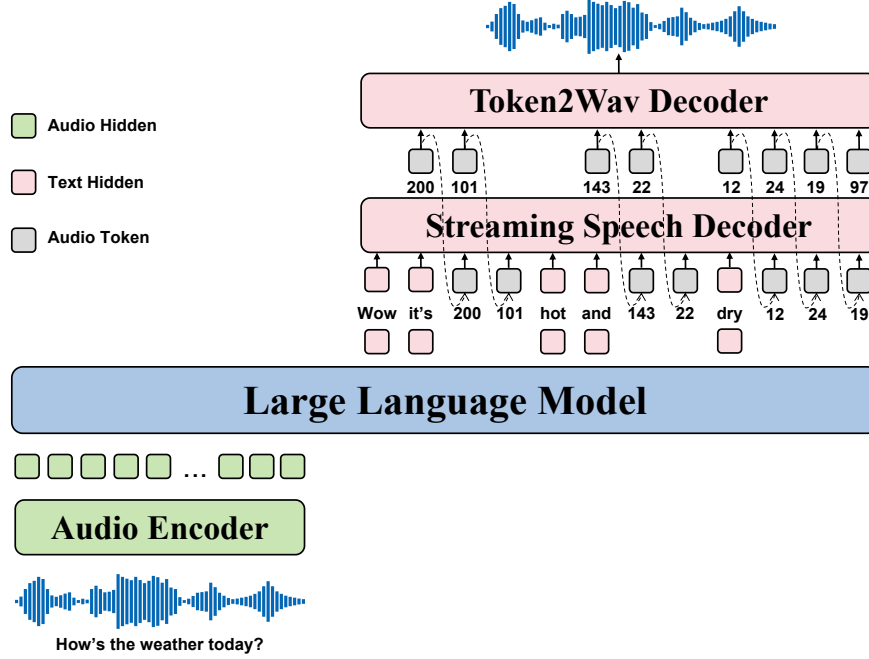
Figure 1: Architecture of our proposed OpenS2S.

## 2 Method

### 2.1 Architecture

The OpenS2S model architecture is shown in Figure 1, comprising four components: an audio encoder, an instruction-following LLM, a streaming speech decoder, and a token2wav decoder. Next, we will describe how to understand continuous speech input and ultimately generate an empathetic speech response.

**Audio Encoder** The Audio Encoder is responsible for transforming this raw audio signal into a more manageable and meaningful representation. To achieve this we use the encoder of Qwen2-Audio [15] to extract features from the audio waveform due to its powerful ability to encode semantic content and paralinguistic information. These features generated by the audio encoder are encoded at a frequency of 25Hz. To further reduce the sequence length, the encoded representations are fed into a speech adapter, which comprises a downsampling module and a feed-forward network. The downsampling module, consisting of two CNN layers, is designed to compress the sequence length by a factor of 4. Finally, the features output by the speech adapter yield continuous encoded representations at 6.25Hz.

**Instruction-Following LLM** The audio embeddings and text embeddings are concatenated to form interleaved input sequences for the large language model. We select Qwen3-8B-Instruct [16] as the LLM, leveraging its robust text processing capabilities.

**Streaming Speech Decoder** To enable streaming speech generation, we adopt a framework inspired by Minmo [17] and LLaMA-Omni2 [18]. The speech response is first converted into discrete tokens using a supervised semantic speech tokenizer. Then, an autoregressive text-to-speech language model is used to generate speech tokens conditioned on the hidden states of the LLM, enabling real-time generation.

The speech tokenizer is implemented by inserting a quantization module into the encoder of Whisper-large-v3 [5], ultimately producing a token sequence at a resolution of 12.5 tokens per second with a vocabulary size of 16,384. We leverage the pretrained speech tokenizer from GLM-4-Voice [12].

3

Once the speech response is tokenized, a decoder-only Transformer models the conditional generation from LLM hidden states to speech tokens. This decoder is initialized from Qwen3-1.8B, with its vocabulary extended to include the 16,384 speech tokens. The input to the streaming speech decoder consists of the final hidden states from the LLM, which are first projected via a linear layer to match the embedding dimension of the speech decoder.

To achieve streaming generation, we interleave the LLM hidden states and generated speech tokens in a predefined ratio: for every $M$ hidden states consumed, $N$ speech tokens are generated (in our implementation, $M = 4$ and $N = 8$). After all hidden states are consumed, the model continues to autoregressively generate the remaining speech tokens until the response is complete. During training, the cross-entropy loss is computed only on the generated speech tokens.

**Token2Wav Decoder**    The speech tokens generated by the streaming speech decoder are subsequently converted into the final speech waveform by the token2wav decoder. This module comprises two key components: a chunk-aware causal flow matching model, which incrementally synthesizes every M speech tokens into mel-spectrograms in a streaming fashion, and a HiFi-GAN vocoder, which converts the mel-spectrograms into the final waveform. Both the flow matching model and the vocoder are adopted from the pretrained components in GLM-4-Voice [12].

## 2.2   Training Strategy

The training of `OpenS2S` consists of three stages: speech understanding pre-training, speech generation pre-training, and empathy speech instruction fine-tuning. In the first two pre-training stages, we utilize open-source ASR and TTS datasets for pre-training to endow the model with robust speech understanding and generation capabilities. In the instruction fine-tuning stage, we construct an empathy speech instruction dataset for fine-tuning, enabling the model to understand the semantic content and paralinguistic cues in speech, and finally generate empathic speech responses.

**Speech Understanding Pretraining (Stage 1 in Figure 2)**    To equip the model with robust speech understanding capabilities, we perform pretraining on large-scale speech-text paired corpora, following the training paradigm of BLSP-Emo [1]. This stage is divided into two phases: semantic alignment and emotional alignment.

In the semantic alignment phase, we adopt the concept of behavioral alignment, requiring the LLM to produce identical continuations when given either speech or its corresponding transcript as input. Specifically, we first prompt the LLM to generate continuations based on text transcripts from an ASR dataset. During training, the model is required to produce the same continuation when conditioned on continuous speech representations, thus aligning the model's behavior across modalities.

In the emotional alignment phase, we leverage an SER dataset where each transcript is annotated with an emotion label. An LLM is prompted to generate emotion-aware continuations based on the transcript and the reference emotion. We then adapt a speech-language model to generate similar continuations directly from the speech input. This step encourages the model to comprehend and reflect both linguistic semantics and paralinguistic emotional cues, producing text responses aligned with those generated by the LLM given identical content and emotion labels.

Throughout this pretraining process, the parameters of the audio encoder and the LLM remain frozen. Only the speech adapter is fine-tuned to facilitate modality bridging.

**Speech Generation Pretraining (Stage 2 in Figure 2)**    As the streaming speech decoder is initialized from Qwen3-1.8B, which is originally designed for text generation, we first perform offline TTS pretraining to enable it to generate discrete speech tokens. Specifically, we expand the decoder's vocabulary to include 16,384 speech tokens. During this phase, the input text is embedded using word embeddings as a prefix, and the target is the sequence of speech tokens extracted from reference audio. This allows the decoder to learn a basic mapping from text tokens to speech tokens, serving as a foundation for downstream speech synthesis.

In the second stage, we further train the speech decoder to integrate with the LLM for streaming generation. Unlike the previous step, the input text is no longer embedded directly into the speech decoder. Instead, it is first processed by the LLM using a structured prompt. The final-layer hidden states corresponding to the response portion are then extracted and interleaved with speech tokens

Figure 2: The training process of OpenS2S.

during training. At this stage, the LLM's parameters are kept frozen, while the linear projection layer and the speech decoder are fine-tuned. This training strategy not only bridges the LLM and the speech decoder but also adapts the offline decoder into a streaming-capable model that supports interleaved token generation.

**Empathetic Speech Instruction Tuning (Stage 3 in Figure 2)** Following the pretraining stages, the model demonstrates general speech understanding capabilities, enabling it to generate empathetic text responses conditioned on both semantic content and emotional cues in speech. However, its speech generation ability remains limited: the model is only able to produce meaningful speech tokens when explicitly prompted with speech synthesis tasks. When handling general-purpose text instructions or directly responding to speech instructions, the speech decoder often fails to generate coherent or meaningful speech outputs.

We attribute this limitation to overfitting during the TTS-based pretraining stage. Specifically, the model learns to rely on a narrow representation subspace defined by the TTS task, resulting in poor generalizability to broader instruction-following scenarios. To address this issue, we introduce an additional instruction tuning stage aimed at enabling robust and flexible speech generation across diverse task types.

In this stage, the speech encoder remains frozen, while all other components, including audio adapter, the LLM, linear projection layer, and speech decoder, are fully fine-tuned. We observe that relying solely on speech-to-speech instruction data is insufficient to generalize speech generation to textual instructions. Therefore, we further incorporate text-to-speech instruction data, allowing the model to handle both speech and text inputs seamlessly. The construction process of these instruction datasets is detailed in Section 3.2.

## 3  Data Collection

### 3.1  Pre-training

**Speech Understanding**   For the semantic alignment stage, we utilize publicly available ASR datasets, and for the emotion alignment stage, we employ standard SER datasets. The ASR datasets include LibriSpeech [19], CommonVoice 13.0 [20], and the GigaSpeech [21] M subset, comprising approximately 1.9 million English (speech, text) pairs. A comparable number of Chinese ASR samples are randomly drawn from WeNetSpeech [22]. The SER datasets consist of IEMOCAP [23], MELD [24], CMU-MOSEI [25], MEAD [26], and ESD [27], collectively covering around 70k utterances in both English and Chinese.

**Speech Generation**   In the first stage of pretraining, aimed at expanding the vocabulary and enhancing the speech decoder's capacity to generate speech tokens, we use 5k hours of English and 5k hours of Chinese speech randomly sampled from the Emilia [28] dataset. For connecting the large language model with the speech decoder and adapting it into an interleaved streaming generation decoder, we further sample 1k hours each of English and Chinese data from the first-stage dataset for training.

### 3.2  Supervised Fine-tuning

Existing open-source speech instruction datasets commonly face three key challenges:

- Limited speaker diversity, which undermines model robustness to varied speech inputs. For example, datasets such as InstructS2S-200K [29] and E-Chat [30] rely on TTS systems to synthesize speech from text, resulting in limited variability in speaker characteristics.
- Neglect of paralinguistic information, with an exclusive focus on semantic content. Although VoiceAssistant-400K [31] introduces speaker diversity via voice cloning, it fails to capture critical paralinguistic cues such as emotion and speaking style.
- Insufficient label granularity, as exemplified by SD-Eval [32], which is divided into subsets that annotate different paralinguistic attributesfeatures in isolation—such as emotion or gender—in isolation. Nowithout any single subset providesoffering joint annotations across multiple paralinguistic dimensions.

To address these limitations, we propose a fully automated framework for constructing an empathetic speech instruction dataset. This framework systematically enhances speech diversity and representativeness across dimensions such as emotion, age, and gender through the integration of heterogeneous data sources. The full pipeline is illustrated in Figure 3 and comprises the following three stages:

**Collection and Manual Annotation of Seed Audio**   We begin by selecting seed audio samples from several publicly available speech emotion recognition datasets. These datasets cover a wide range of emotions and speaker demographics, including children, adults, and the elderly, ensuring strong representativeness and diversity. Each selected sample is manually annotated with its transcribed text, speaker gender, age, and emotional label. This results in 1,000 English and 1,000 Chinese seed audio samples with rich multi-dimensional annotations.

**Generation of Speech Instructions**   A straightforward solution to generating speech instructions is to convert existing text instruction datasets into speech. However, such datasets present two challenges: (1) tasks involving math or programming are often unsuitable for conversational speech scenarios; and (2) most instruction data neglect paralinguistic factors. Inspired by the Self-Instruct paradigm, we leverage Qwen3-32B-Instruct to automatically generate task instructions that are
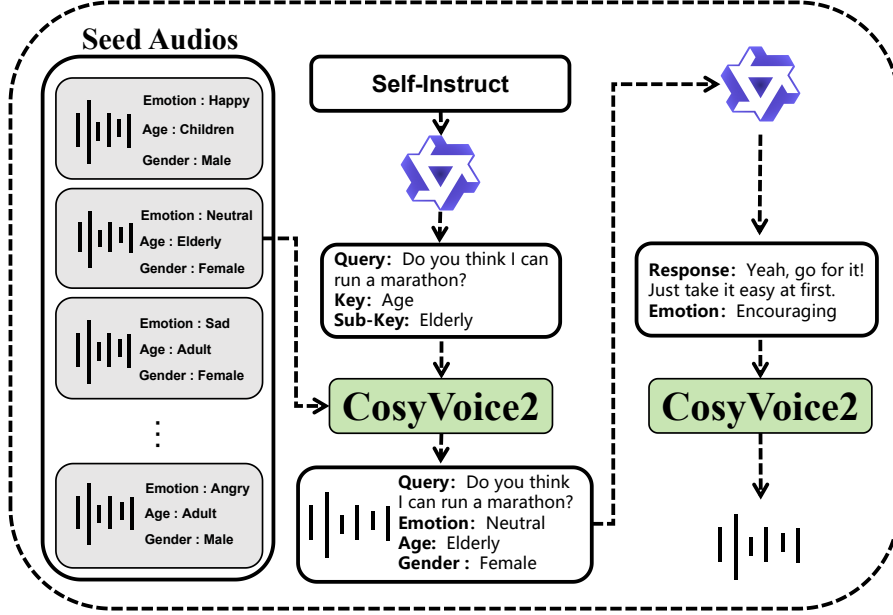
Figure 3: Automatic workflow for constructing empathetic speech instruction dataset.

sensitive to paralinguistic cues. For example, the model might generate the instruction "Do you think I can run a marathon?" and tag it as age-sensitive, suggesting it be read in an elderly voice. We then randomly select a seed audio whose emotion label is "elderly" as an audio prompt and use CosyVoice2 [33] for voice cloning, preserving the emotion and gender attributes of the selected seed. This process yields 50k English and 50k Chinese speech instructions with diverse paralinguistic characteristics.

**Generation of Speech Response** For each speech instruction, we annotate its transcribed text, emotion, age, and gender based on the matched seed audio. Both the semantic content and paralinguistic labels are input into Qwen3-32B-Instruct with "thinking mode" enabled. Inspired by LLaMA-Omni, we prompt the model to generate concise, dialogue-appropriate, and empathetic text responses. We then prompt the same model to infer the appropriate emotional tone for delivering the response, conditioned on the original instruction, the response content, and the paralinguistic features of the input speech. Finally, we use a consistent reference voice as a prompt and control CosyVoice2 via instruction prompts to synthesize emotionally expressive speech responses.

Through this systematic construction process, we obtain an empathetic speech instruction dataset characterized by multi-dimensional emotion annotation, expressive emotional delivery, and diverse speaker profiles. The statistics of the constructed data are presented in Figure 4. In total, we construct 50k English and 50k Chinese speech-to-speech empathetic samples. The input speech includes three types of paralinguistic information tags (i.e. emotion, gender, age)), and the output speech is fixed as a young female voice responding with different emotions.

To retain the model's general instruction-following capabilities beyond empathetic dialogue, we further incorporate general-purpose data. Specifically, we extract 50k English instructions from Instruct200K [29], translate them into Chinese using Qwen3-32B-Instruct, and select seed audio with neutral emotional labels to convert these instructions into speech. Applying the same speech response generation process as above, we obtain an additional 100k bilingual speech-to-speech instruction pairs.

Lastly, we observe that training solely on speech-to-speech data can impair the model's ability to respond to text inputs: while the language model remains capable of generating reasonable text, the speech decoder fails to produce valid speech tokens when conditioned on text instruction. To mitigate this, we extract text-to-speech instruction samples from the general speech-to-speech dataset, ensuring the model can jointly handle both text and speech inputs during inference.
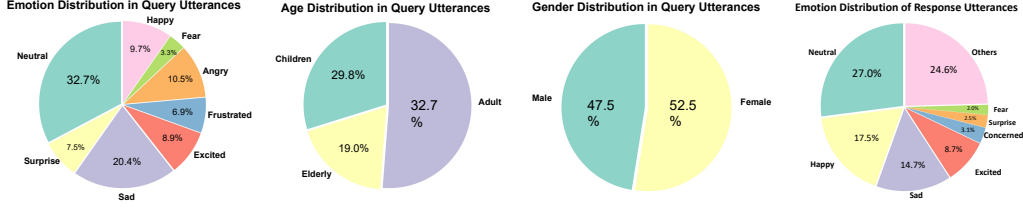
Figure 4: Distributions of emotion, age, and gender in query utterances, and emotion distribution in response utterances.

Table 2: Performance of `OpenS2S` and baseline models on the benchmarks of speech-to-text chat.

| Model | VoiceBench | | | | URO-Bench | |
|---|---|---|---|---|---|---|
| | alpaca | common | ifeval | wildvoice | underemo-en | underemo-zh |
| Qwen-2-Audio [15] | 3.74 | 3.43 | 26.33 | 3.01 | 35.38 | 69.62 |
| GLM-4-Voice [12] | 3.97 | 3.42 | 25.92 | 3.18 | 52.41 | 74.51 |
| Kimi-Audio [13] | 4.46 | 3.97 | 61.10 | 4.20 | 59.22 | 76.96 |
| LLaMA-Omni2[2] [18] | 3.96 | 3.46 | 17.36 | 3.07 | 39.46 | 63.79 |
| OpenS2S | 4.09 | 3.65 | 42.89 | 3.66 | 46.90 | 67.68 |

# 4 Evaluation

## 4.1 Speech-to-Text Chat

We evaluate the ability of `OpenS2S` to engage in speech-to-text conversations based on audio input using two benchmarks:

- **VoiceBench** [34] is a benchmark designed for the multi-faceted evaluation of LLM-based voice assistants. To assess the crucial capability of instruction-following, we utilize the alpacaeval, commoneval, wildvoice, and ifeval subsets. These were specifically chosen to test the model's ability to comprehend and accurately execute diverse spoken commands.

- **URO-Bench** [35] is an end-to-end benchmark for spoken dialogue models that assesses understanding, reasoning, and oral conversation skills, including paralinguistic cues. To evaluate the model's capacity for empathy, we employ its UnderEmotion-en and UnderEmotion-zh subsets. These are designed to measure the model's ability to perceive the user's emotional state and generate affectively appropriate responses in both English and Chinese.

These two benchmarks evaluate from multiple perspectives to ensure comprehensive assessment: VoiceBench assesses the model's instruction-following capability, while URO-Bench evaluates its comprehension and response to paralinguistic emotional cues. We also evaluate several mainstream open-source speech large models with comparable parameter sizes to `OpenS2S` for comparison. The results are presented in Table 2.

The results in Table 2 show that `OpenS2S` demonstrates competitive performance across the four subsets of VoiceBench. Its scores rank second only to Kimi-Audio, which is trained on substantially more data[3], and outperform all other models. These findings indicate that `OpenS2S` possesses strong capabilities in spoken dialogue and can effectively handle user voice command inputs.

In addition, the results on the URO-Bench subsets show that `OpenS2S` achieves scores close to those of state-of-the-art models in empathy evaluation, despite being trained on significantly less data. This not only confirms the solid empathetic interaction capabilities of `OpenS2S`, but also highlights the high quality of the data generated by the proposed empathetic speech dialogue data generation method.

---

[2]We compare with LLaMA-Omni2-7B-Bilingual for its comparable parameter size and bilingual capability.
[3]For example, Kimi-Audio employs more 13 million hours of audio data for pre-training.

### 4.2 Speech-to-Speech Chat

Finally, we assess the end-to-end speech conversation capabilities of `OpenS2S` based on qualitative analysis. Visit `https://casia-lm.github.io/OpenS2S` for demos.

## 5 Related Work

### 5.1 Speech Language Models

With the rapid advancement of large language models (LLMs), there is growing interest in extending their capabilities to spoken language, giving rise to Speech Language Models (SpeechLMs) that can understand and/or generate speech [36, 37]. One approach directly adapts LLMs for *end-to-end speech modeling* by converting speech into discrete tokens and expanding the vocabulary, as seen in SpeechGPT [38], AudioPaLM [39], and TWIST [40]. Recent models like Spirit-LM [41] and GLM-4-Voice [12] leverage interleaved speech-text training, while others such as Moshi [11] and LSLM [42] enable spoken dialogue. In contrast, modular SpeechLMs *connect LLMs with external speech modules*. Early works [43–48, 15, 49] focused on speech understanding by connecting pretrained speech encoders to LLMs, but did not support speech generation. In contrast, more recent models such as LLaMA-Omni [29, 18], Freeze-Omni [50], and OpenOmni [51] overcome this limitation by attaching speech decoders to LLM outputs. Mini-Omni [31] and SLAM-Omni [52] go further with parallel decoding. Minmo [17] and LLaMA-Omni2 [18] incorporates a streaming speech decoder through interleaved text-speech generation.

### 5.2 Empathetic Conversations Across Modalities

Empathetic conversation modeling [3, 53] has been studied across *text-to-text, speech-to-text, and speech-to-speech* settings, aiming to equip LLMs with emotional understanding and supportive responses [54]. In text-based interactions, early work focused on architecture modifications [55], while recent approaches like SoulChat [56] and Chain of Empathy prompting [57] enhance empathy through fine-tuning or step-by-step reasoning without extra data. For speech-to-text interaction, E-chat [30] introduced an emotion-aware speech instruction dataset to enhance LLMs' understanding of emotional speech. BLSP-Emo [1] proposed an end-to-end model that aligns speech semantics and emotions through two-stage pretraining using ASR and SER datasets. Moving toward speech-to-speech empathy, Spoken-GPT [58] adopts a cascaded framework that listens and responds with expressive, emotionally attuned speech, paving the way for fully empathetic voice agents. Advanced commercial models such as GPT-4o [14], Doubao, Kimi-Audio [13], and Step-Audio [59] push the boundaries of empathetic interaction by incorporating paralinguistic cues to better perceive and respond to users' emotional states. These models integrate speech understanding and generation in real time, enabling more natural and emotionally aware human-computer interactions. However, our `OpenS2S` is the first to release all the resources including model weights, training data and training codes, in order to boost the research in the community.

## 6 Conclusion

This report presents `OpenS2S`, a fully open-source, end-to-end LSLM specifically designed for empathetic speech interactions. `OpenS2S` distinguishes itself with an efficient streaming interleaved decoding architecture, enabling low-latency response generation, and an innovative automated data construction pipeline. This pipeline cost-effectively synthesizes diverse, high-quality empathetic speech dialogues by leveraging large language models and controllable text-to-speech systems. As a result, `OpenS2S` achieves competitive performance in empathetic interactions while requiring substantially less data and computational resources compared to current resource-intensive pre-training methods. We release the complete `OpenS2S` framework, including the dataset, model weights, and training codes, to empower the broader research community and accelerate innovation in empathetic speech systems.

# References

[1] Chen Wang, Minpeng Liao, Zhongqiang Huang, Junhong Wu, Chengqing Zong, and Jiajun Zhang. Blsp-emo: Towards empathetic large speech-language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 19186–19199, 2024.

[2] Sylvia A Morelli, Matthew D Lieberman, and Jamil Zaki. The emerging study of positive empathy. *Social and Personality Psychology Compass*, 9(2):57–68, 2015.

[3] Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. Towards empathetic open-domain conversation models: A new benchmark and dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, 2019.

[4] Anjuli Kannan, Arindrima Datta, Tara N Sainath, Eugene Weinstein, Bhuvana Ramabhadran, Yonghui Wu, Ankur Bapna, Zhifeng Chen, and Seungji Lee. Large-scale multilingual speech recognition with a streaming end-to-end model. *arXiv preprint arXiv:1909.05330*, 2019.

[5] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR, 2023.

[6] Hamza Kheddar, Mustapha Hemis, and Yassine Himeur. Automatic speech recognition using advanced deep learning approaches: A survey. *Information Fusion*, page 102422, 2024.

[7] Hongkun Hao, Long Zhou, Shujie Liu, Jinyu Li, Shujie Hu, Rui Wang, and Furu Wei. Boosting large language model for speech synthesis: An empirical study. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2025.

[8] Yinghao Aaron Li, Cong Han, Vinay Raghavan, Gavin Mischler, and Nima Mesgarani. Styletts 2: Towards human-level text-to-speech through style diffusion and adversarial training with large speech language models. *Advances in Neural Information Processing Systems*, 36:19594–19621, 2023.

[9] Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Mark Duppenthaler, Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim, et al. Seamless: Multilingual expressive and streaming speech translation. *arXiv preprint arXiv:2312.05187*, 2023.

[10] Zhichao Huang, Rong Ye, Tom Ko, Qianqian Dong, Shanbo Cheng, Mingxuan Wang, and Hang Li. Speech translation with large language models: An industrial practice. *arXiv preprint arXiv:2312.13585*, 2023.

[11] Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard Grave, and Neil Zeghidour. Moshi: a speech-text foundation model for real-time dialogue. *arXiv preprint arXiv:2410.00037*, 2024.

[12] Aohan Zeng, Zhengxiao Du, Mingdao Liu, Kedong Wang, Shengmin Jiang, Lei Zhao, Yuxiao Dong, and Jie Tang. Glm-4-voice: Towards intelligent and human-like end-to-end spoken chatbot. *arXiv preprint arXiv:2412.02612*, 2024.

[13] Ding Ding, Zeqian Ju, Yichong Leng, Songxiang Liu, Tong Liu, Zeyu Shang, Kai Shen, Wei Song, Xu Tan, Heyi Tang, et al. Kimi-audio technical report. *arXiv preprint arXiv:2504.18425*, 2025.

[14] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.

[15] Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, et al. Qwen2-audio technical report. *arXiv preprint arXiv:2407.10759*, 2024.

[16] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

[17] Qian Chen, Yafeng Chen, Yanni Chen, Mengzhe Chen, Yingda Chen, Chong Deng, Zhihao Du, Ruize Gao, Changfeng Gao, Zhifu Gao, et al. Minmo: A multimodal large language model for seamless voice interaction. *arXiv preprint arXiv:2501.06282*, 2025.

[18] Qingkai Fang, Yan Zhou, Shoutao Guo, Shaolei Zhang, and Yang Feng. Llama-omni2: Llm-based real-time spoken chatbot with autoregressive streaming speech synthesis. *arXiv preprint arXiv:2505.02625*, 2025.

[19] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: an asr corpus based on public domain audio books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5206–5210. IEEE, 2015.

[20] Rosana Ardila, Megan Branson, Kelly Davis, Michael Kohler, Josh Meyer, Michael Henretty, Reuben Morais, Lindsay Saunders, Francis Tyers, and Gregor Weber. Common voice: A massively-multilingual speech corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4218–4222, 2020.

[21] Guoguo Chen, Shuzhou Chai, Guanbo Wang, Jiayu Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel Povey, Jan Trmal, Junbo Zhang, et al. Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio. *arXiv preprint arXiv:2106.06909*, 2021.

[22] Binbin Zhang, Hang Lv, Pengcheng Guo, Qijie Shao, Chao Yang, Lei Xie, Xin Xu, Hui Bu, Xiaoyu Chen, Chenchen Zeng, et al. Wenetspeech: A 10000+ hours multi-domain mandarin corpus for speech recognition. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6182–6186. IEEE, 2022.

[23] Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N Chang, Sungbok Lee, and Shrikanth S Narayanan. Iemocap: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42:335–359, 2008.

[24] Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. Meld: A multimodal multi-party dataset for emotion recognition in conversations. *arXiv preprint arXiv:1810.02508*, 2018.

[25] AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable dynamic fusion graph. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2236–2246, 2018.

[26] Kaisiyuan Wang, Qianyi Wu, Linsen Song, Zhuoqian Yang, Wayne Wu, Chen Qian, Ran He, Yu Qiao, and Chen Change Loy. Mead: A large-scale audio-visual dataset for emotional talking-face generation. In *European Conference on Computer Vision*, pages 700–717. Springer, 2020.

[27] Kun Zhou, Berrak Sisman, Rui Liu, and Haizhou Li. Emotional voice conversion: Theory, databases and esd. *Speech Communication*, 137:1–18, 2022.

[28] Haorui He, Zengqiang Shang, Chaoren Wang, Xuyuan Li, Yicheng Gu, Hua Hua, Liwei Liu, Chen Yang, Jiaqi Li, Peiyang Shi, et al. Emilia: An extensive, multilingual, and diverse speech dataset for large-scale speech generation. In *2024 IEEE Spoken Language Technology Workshop (SLT)*, pages 885–890. IEEE, 2024.

[29] Qingkai Fang, Shoutao Guo, Yan Zhou, Zhengrui Ma, Shaolei Zhang, and Yang Feng. Llama-omni: Seamless speech interaction with large language models. *arXiv preprint arXiv:2409.06666*, 2024.

[30] Hongfei Xue, Yuhao Liang, Bingshen Mu, Shiliang Zhang, Mengzhe Chen, Qian Chen, and Lei Xie. E-chat: Emotion-sensitive spoken dialogue system with large language models. In *2024 IEEE 14th International Symposium on Chinese Spoken Language Processing (ISCSLP)*, pages 586–590. IEEE, 2024.

[31] Zhifei Xie and Changqiao Wu. Mini-omni: Language models can hear, talk while thinking in streaming. *arXiv preprint arXiv:2408.16725*, 2024.

[32] 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. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

[33] Zhihao Du, Yuxuan Wang, Qian Chen, Xian Shi, Xiang Lv, Tianyu Zhao, Zhifu Gao, Yexin Yang, Changfeng Gao, Hui Wang, et al. Cosyvoice 2: Scalable streaming speech synthesis with large language models. *arXiv preprint arXiv:2412.10117*, 2024.

[34] Yiming Chen, Xianghu Yue, Chen Zhang, Xiaoxue Gao, Robby T. Tan, and Haizhou Li. Voicebench: Benchmarking llm-based voice assistants. *arXiv preprint arXiv:2410.17196*, 2024.

[35] Ruiqi Yan, Xiquan Li, Wenxi Chen, Zhikang Niu, Chen Yang, Ziyang Ma, Kai Yu, and Xie Chen. Uro-bench: A comprehensive benchmark for end-to-end spoken dialogue models. *arXiv preprint arXiv:2502.17810*, 2025.

[36] Wenqian Cui, Dianzhi Yu, Xiaoqi Jiao, Ziqiao Meng, Guangyan Zhang, Qichao Wang, Yiwen Guo, and Irwin King. Recent advances in speech language models: A survey. *arXiv preprint arXiv:2410.03751*, 2024.

[37] Shengpeng Ji, Yifu Chen, Minghui Fang, Jialong Zuo, Jingyu Lu, Hanting Wang, Ziyue Jiang, Long Zhou, Shujie Liu, Xize Cheng, et al. Wavchat: A survey of spoken dialogue models. *arXiv preprint arXiv:2411.13577*, 2024.

[38] Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15757–15773, 2023.

[39] Paul K Rubenstein, Chulayuth Asawaroengchai, Duc Dung Nguyen, Ankur Bapna, Zalán Borsos, Félix de Chaumont Quitry, Peter Chen, Dalia El Badawy, Wei Han, Eugene Kharitonov, et al. Audiopalm: A large language model that can speak and listen. *arXiv preprint arXiv:2306.12925*, 2023.

[40] Michael Hassid, Tal Remez, Tu Anh Nguyen, Itai Gat, Alexis Conneau, Felix Kreuk, Jade Copet, Alexandre Defossez, Gabriel Synnaeve, Emmanuel Dupoux, et al. Textually pretrained speech language models. *Advances in Neural Information Processing Systems*, 36:63483–63501, 2023.

[41] Tu Anh Nguyen, Benjamin Muller, Bokai Yu, Marta R Costa-Jussa, Maha Elbayad, Sravya Popuri, Christophe Ropers, Paul-Ambroise Duquenne, Robin Algayres, Ruslan Mavlyutov, et al. Spirit-lm: Interleaved spoken and written language model. *Transactions of the Association for Computational Linguistics*, 13:30–52, 2025.

[42] Ziyang Ma, Yakun Song, Chenpeng Du, Jian Cong, Zhuo Chen, Yuping Wang, Yuxuan Wang, and Xie Chen. Language model can listen while speaking. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 24831–24839, 2025.

[43] Yu Shu, Siwei Dong, Guangyao Chen, Wenhao Huang, Ruihua Zhang, Daochen Shi, Qiqi Xiang, and Yemin Shi. Llasm: Large language and speech model. *arXiv preprint arXiv:2308.15930*, 2023.

[44] Chen Wang, Minpeng Liao, Zhongqiang Huang, Jinliang Lu, Junhong Wu, Yuchen Liu, Chengqing Zong, and Jiajun Zhang. Blsp: Bootstrapping language-speech pre-training via behavior alignment of continuation writing. *arXiv preprint arXiv:2309.00916*, 2023.

[45] Chen Wang, Minpeng Liao, Zhongqiang Huang, and Jiajun Zhang. Blsp-kd: Bootstrapping language-speech pre-training via knowledge distillation. *arXiv preprint arXiv:2405.19041*, 2024.

[46] Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*, 2023.

[47] Eliya Nachmani, Alon Levkovitch, Roy Hirsch, Julian Salazar, Chulayuth Asawaroengchai, Soroosh Mariooryad, Ehud Rivlin, RJ Skerry-Ryan, and Michelle Tadmor Ramanovich. Spoken question answering and speech continuation using spectrogram-powered llm. In *The Twelfth International Conference on Learning Representations*.

[48] Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, MA Zejun, and Chao Zhang. Salmonn: Towards generic hearing abilities for large language models. In *The Twelfth International Conference on Learning Representations*.

[49] Shujie Hu, Long Zhou, Shujie Liu, Sanyuan Chen, Lingwei Meng, Hongkun Hao, Jing Pan, Xunying Liu, Jinyu Li, Sunit Sivasankaran, et al. Wavllm: Towards robust and adaptive speech large language model. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4552–4572, 2024.

[50] Xiong Wang, Yangze Li, Chaoyou Fu, Yunhang Shen, Lei Xie, Ke Li, Xing Sun, and Long Ma. Freeze-omni: A smart and low latency speech-to-speech dialogue model with frozen llm. *arXiv preprint arXiv:2411.00774*, 2024.

[51] Run Luo, Ting-En Lin, Haonan Zhang, Yuchuan Wu, Xiong Liu, Min Yang, Yongbin Li, Longze Chen, Jiaming Li, Lei Zhang, et al. Openomni: Large language models pivot zero-shot omnimodal alignment across language with real-time self-aware emotional speech synthesis. *arXiv preprint arXiv:2501.04561*, 2025.

[52] Wenxi Chen, Ziyang Ma, Ruiqi Yan, Yuzhe Liang, Xiquan Li, Ruiyang Xu, Zhikang Niu, Yanqiao Zhu, Yifan Yang, Zhanxun Liu, et al. Slam-omni: Timbre-controllable voice interaction system with single-stage training. *arXiv preprint arXiv:2412.15649*, 2024.

[53] Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie Huang. Towards emotional support dialog systems. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3469–3483, 2021.

[54] Brant R Burleson. Emotional support skills. In *Handbook of communication and social interaction skills*, pages 569–612. Routledge, 2003.

[55] Raman Goel, Seba Susan, Sachin Vashisht, and Armaan Dhanda. Emotion-aware transformer encoder for empathetic dialogue generation. In *2021 9th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos (ACIIW)*, pages 1–6. IEEE, 2021.

[56] Yirong Chen, Xiaofen Xing, Jingkai Lin, Huimin Zheng, Zhenyu Wang, Qi Liu, and Xiangmin Xu. Soulchat: Improving llms' empathy, listening, and comfort abilities through fine-tuning with multi-turn empathy conversations. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1170–1183, 2023.

[57] Yoon Kyung Lee, Inju Lee, Minjung Shin, Seoyeon Bae, and Sowon Hahn. Chain of empathy: Enhancing empathetic response of large language models based on psychotherapy models. *arXiv preprint arXiv:2311.04915*, 2023.

[58] Dong Zhang, Xin Zhang, Jun Zhan, Shimin Li, Yaqian Zhou, and Xipeng Qiu. Speechgpt-gen: Scaling chain-of-information speech generation. *arXiv preprint arXiv:2401.13527*, 2024.

[59] Ailin Huang, Boyong Wu, Bruce Wang, Chao Yan, Chen Hu, Chengli Feng, Fei Tian, Feiyu Shen, Jingbei Li, Mingrui Chen, et al. Step-audio: Unified understanding and generation in intelligent speech interaction. *arXiv preprint arXiv:2502.11946*, 2025.