

BrainSpeech: Parallel Semantic–Acoustic Generation for SpeechLMs

Anonymous ACL submission

Abstract

Recent advances in Speech Language Models (SpeechLMs) have enabled end-to-end spoken dialogue systems. However, most existing SpeechLMs rely on a single-stream or sequential modeling paradigm, in which semantic reasoning and acoustic generation are either flattened into one stream or processed in a fixed Thinker–Talker order. Both paradigms constrain the interaction between semantic reasoning and acoustic generation, limiting the naturalness and expressiveness of spoken dialogue. To address this, we propose BrainSpeech, a novel SpeechLM framework that models semantic reasoning and acoustic generation as two parallel yet interacting streams within a unified large language model. To enable effective cross-stream interaction, we introduce a dedicated attention mechanism that explicitly exchanges semantic and acoustic information during generation. Building on this dual-stream formulation, we further develop a three-stage training strategy that preserves the reasoning capability of the underlying language model while reducing the reliance on large-scale spoken dialogue corpora. In addition, we design a streaming decoding strategy that supports real-time generation of continuous, high-fidelity speech. Experiments on multiple spoken dialogue benchmarks demonstrate that BrainSpeech produces more natural and expressive speech and achieves superior speech-to-speech performance compared to larger open-source SpeechLMs¹.

1 Introduction

Spoken dialogue systems (Long et al., 2025) play a crucial role in natural human–computer interaction by lowering interaction barriers and enabling more intuitive and seamless user experiences. Traditional speech interaction pipelines (Wang

et al., 2024a) typically consist of three dedicated components: Automatic Speech Recognition (ASR) (Saini and Kau, 2013) for transcribing speech, a large language model (LLM) (Tsai et al., 2023) for generating textual responses, and Text-to-Speech (TTS) (Liu et al., 2025) for synthesizing speech outputs. However, such cascaded architectures suffer from accumulated interaction latency and limited modeling of emotional prosody. To overcome these limitations, recent studies have increasingly explored end-to-end Speech Language Model (SpeechLM)-based frameworks that directly map speech inputs to speech responses (Xie and Wu, 2024a; Chen et al., 2025b; Fang et al., 2024). In these systems, speech inputs are first encoded by an audio encoder and aligned to the LLM’s textual representation space via an adaptor, after which speech responses are generated using either discrete audio tokens (Zeng et al., 2024) or continuous hidden states (Wang et al., 2024b) produced by the LLM.

Existing end-to-end SpeechLMs can be categorized into sequential-stream and single-stream-based systems depending on the manner in which understanding and generation are carried out. Sequential-stream-based systems (Yang et al., 2024; Chen et al., 2025a; Long et al., 2025; Wang et al., 2024b; Fang et al., 2025) perform text-based semantic reasoning before generating speech from text tokens or continuous semantic representations. These systems better preserve the comprehension capabilities of LLMs but tend to lose paralinguistic information such as emotional and prosodic cues. Single-stream-based systems (Zeng et al., 2024; Chen et al., 2025b; Xie and Wu, 2024b; KimiTeam et al., 2025; Wu et al., 2025) employ a unified LLM to simultaneously perform thinking and generation, generate text and audio tokens in a parallel or interleaved manner. Although these systems are advantageous for modeling paralinguistic information,

¹Audio samples are available at <https://brainspeech.github.io/>

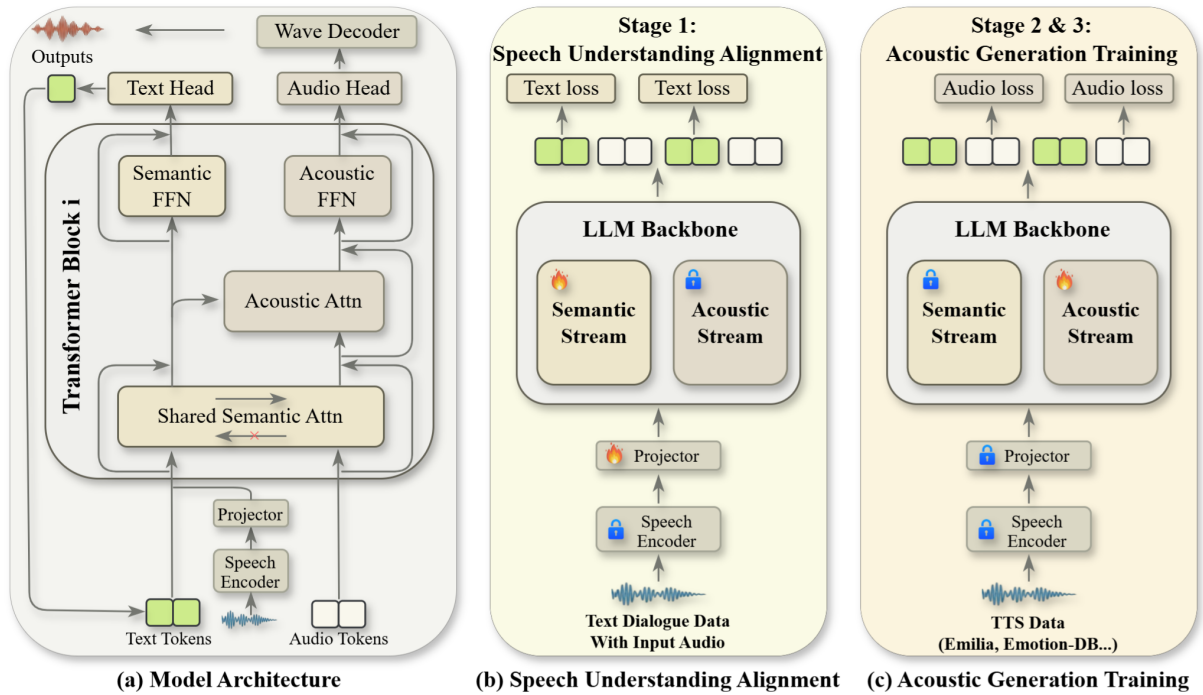


Figure 1: The overall framework and training paradigm of BrainSpeech. (a) Model Architecture. Text and audio tokens are processed by the Semantic Stream (Shared Semantic Attn + Semantic FFN) and Acoustic Stream (Shared Semantic Attn + Acoustic Attn + Acoustic FFN), respectively. (b) Stage 1: Speech Understanding Alignment. Trains the Semantic Stream and Projector to enable speech understanding. (c) Stages 2 & 3: Acoustic Generative Training. Trains the Acoustic Stream on TTS data to enable speech generation.

3.1 Model Architecture

As illustrated in Figure 1, the architecture of BrainSpeech is a dual-stream architecture LLM that maintains semantic and acoustic processing streams within each block. To enable speech understanding, speech queries are first processed by the pre-trained speech encoder (Radford et al., 2023) and mapped to the LLM’s space through a trainable linear projector. Then, the hidden states of input speech and text tokens are processed sequentially through semantic attention and Feed-Forward Networks (FFN) inherited from LLMs, focusing exclusively on semantic generation. In contrast, the output audio tokens traverse a specialized path for acoustic generation. They first utilize the shared semantic attention to extract high-level semantic guidance, followed by a dedicated acoustic attention to model acoustic information (e.g. prosody and pitch) and finally acoustic FFN to integrate and refine the modeled acoustic information for coherent acoustic representation synthesis. This design ensures the effective integration of semantic and acoustic information. Finally, the generated discrete audio tokens are converted into continuous speech waveforms using the decoder of XCodec (Ye et al., 2025).

3.2 Sequence Format and Terminal Control

In natural human communication, textual responses are typically accompanied by corresponding speech signals to enhance interactive naturalness. Therefore, BrainSpeech generates responses using an interleaved sequence format:

$$\mathbf{T}_{s1}, \langle b \rangle, \mathbf{A}_{s1}, \langle e \rangle, \dots, \mathbf{T}_{sl}, \langle lb \rangle, \mathbf{A}_{sl}, \langle e \rangle \quad (1)$$

To support efficient streaming generation, BrainSpeech imposes a fixed token budget (10:40) for text and audio tokens within intermediate segments. However, a critical challenge is to prevent additional audio token generation in intermediate segments, as this would disrupt the interleaved structure and cause a divergence from the training distribution. We address this by introducing a special token $\langle lb \rangle$ (last begin) to explicitly identify the final segment. This control mechanism strictly constrains generation to a fixed token count for intermediate segments, effectively deferring the generation of remaining audio tokens until the $\langle lb \rangle$ token is predicted, preserving structural integrity during inference.

3.3 Asymmetric Causal Masking

Given the interleaved format, BrainSpeech employs an asymmetric attention mask to disentangle the generation dependencies. Since *what to say* (text) should remain independent of *how to say* (audio), BrainSpeech masks text tokens from attending to audio tokens. This unidirectional constraint preserves the purity of semantic generation and prevents contamination from the update of acoustic features. Formally, let $\mathbf{H}^l = [\mathbf{h}_1^l, \dots, \mathbf{h}_N^l]$ denote the hidden states of layer l . For the text token at index i , the update process is constrained to:

$$\begin{aligned} \tilde{\mathbf{h}}_i^l &= \text{SemanticAttn}(\mathbf{h}_i^l, \{\mathbf{h}_j^l : j \in \mathbf{T}, j \leq i\}) \\ \mathbf{h}_i^{l+1} &= \text{SemanticFFN}(\tilde{\mathbf{h}}_i^l) \end{aligned} \quad (2)$$

Audio tokens utilize a standard causal mask to attend to all preceding text and audio context, thereby ensuring the acoustic features are semantically aligned with text while maintaining prosodic coherence.

3.4 Training Paradigm

For the training paradigm, we employ a progressive three-stage training pipeline. This strategy eliminates the need for large-scale spoken dialogue data by decoupling the optimization of understanding and generation capabilities.

Stage 1: Speech Understanding Alignment. The primary goal is to equip LLMs with the capabilities of speech understanding. We denote the parameters of the projector as $\theta_{\mathcal{P}}$, the semantic stream as $\theta_{\mathcal{S}}$, and the acoustic stream as $\theta_{\mathcal{A}}$. In this stage, we freeze the speech encoder and maximize the likelihood of generating the correct text response \mathbf{T} given the input speech query \mathbf{S} . The optimization objective is:

$$\mathcal{L}_{\text{Stage 1}} = - \sum_{i=1}^{|\mathbf{T}|} \log P(\mathbf{T}_i | \mathbf{T}_{<i}, \mathbf{S}; \theta_{\mathcal{P}}, \theta_{\mathcal{S}}) \quad (3)$$

where $\theta_{\mathcal{A}}$ is not involved, and BrainSpeech is optimized to bridge the modality gap in this stage.

Stage 2: Acoustic Generation Pre-training. This stage trains the acoustic stream to generate audio tokens conditioned on the text response \mathbf{T} . Specifically, we freeze $\theta_{\mathcal{S}}$ to preserve the capabilities of understanding and only update $\theta_{\mathcal{A}}$. The training objective is to maximize the conditional log-likelihood of the target audio sequence

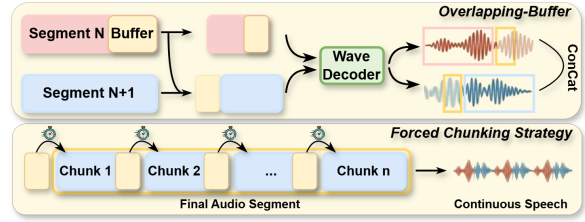


Figure 2: Inference of BrainSpeech.

Given \mathbf{T} . We minimize the following negative log-likelihood loss:

$$\mathcal{L}_{\text{Stage 2\&3}} = - \sum_{t=1}^{|\mathbf{A}|} \log P(\mathbf{A}_t | \mathbf{A}_{<t}, \mathbf{T}; \theta_{\mathcal{A}}) \quad (4)$$

This stage establishes a robust mapping from semantic content to acoustic features.

Stage 3: Expressive Acoustic Generation Fine-tuning. To enhance both acoustic fidelity and emotional expressiveness, we fine-tune $\theta_{\mathcal{A}}$ using a composite dataset. Similar to Stage 2, $\theta_{\mathcal{S}}$ remains frozen. This stage enables BrainSpeech to produce high-quality speech, effectively capturing context-aware prosody and emotion inferred from the text.

3.5 Inference of BrainSpeech

Generating long-form audio in interleaved segments would introduce acoustic artifacts. As shown in Figure 2, we develop specific decoding strategies to ensure high-fidelity speech while enabling continuous real-time streaming.

Overlapping-Buffer: To address boundary discontinuities caused by independent decoding in a segment, we introduce overlapping-buffer. Specifically, we concatenate the trailing tokens of the preceding segment to the current segment, ensuring continuity between segments. We achieve seamless transitions by splicing waveforms at the midpoint of the buffer, preserving the initial half from the prior segment while utilizing the second half from the current segment. Finally, the buffer is refreshed with the current segment’s trailing tokens for the next iteration.

Forced-Chunking: The final audio segment frequently exhibits significant length variation. To ensure stable and real-time decoding, we employ forced-chunking. Tokens are accumulated and processed in chunks aligned with the intermediate segment size, iteratively applying overlapping-buffer to guarantee robust generation independent of total duration.

4 Data Construction

To train BrainSpeech effectively across understanding and generation tasks without relying on large-scale spoken dialogue data, we construct specialized datasets for each training stage. Detailed formatting examples for each stage are provided in Appendix A.

4.1 Data of Speech Understanding Alignment

To align the speech encoder with LLMs, we derive data from high-quality text-only instruction datasets (e.g., InstructS2S (Fang et al., 2024)). We employ commercial TTS models (Higgs (Boson AI, 2025) and Minimax-Speech-02) to synthesize questions into speech. Then, the original text query is replaced by the continuous speech embeddings extracted by the speech encoder. This effectively enables the model to comprehend speech by leveraging its established textual understanding capabilities.

4.2 Data of Acoustic Generation Pre-training

For acoustic generation pre-training, we utilize Emilia (He et al., 2024), a TTS dataset with 2w hours, as training data. To follow the dialogue format required by our architecture, we implement a shuffle-and-pair strategy. Specifically, we randomly sample two distinct audio-text pairs: one serves as the audio "query" and the other as the text and audio "response". Although the context is intentionally randomized, this design forces the model to focus purely on the mapping from the current response text to its corresponding audio tokens, ensuring the stability of audio generation regardless of text context.

4.3 Data of Expressive Acoustic Generation Fine-tuning

To achieve high-fidelity and expressive synthesis, we create a composite dataset, aiming to equip BrainSpeech with both superior acoustic quality and rich emotional expressiveness. We organize the data into two categories:

- *High Quality Dialogue Data:* We leverage the data utilized in stage 1. This process ensures BrainSpeech masters the generation of speech with superior acoustic quality.
- *Expressive Speech Corpora:* Emotion-DB (Yang et al., 2025), a corpus of high-fidelity human speech characterized by rich

emotional dynamics, is used. Therefore, θ_A can capture intricate paralinguistic nuances, thereby enriching the generated responses with diverse and vivid emotional tones.

5 Experimental Setup

5.1 Model Configuration

LLM Backbone and Speech Encoder: We adopt Llama-3.2-3B-Instruct as our base LLM. For the speech encoder, we utilize Whisper-large-v3. The speech encoder are frozen throughout the training process to preserve the capabilities of robust acoustic feature extraction. To bridge the modality gap, we employ a trainable projector consisting of a two-layer MLP.

Initialization of Dual-Stream Architecture: The parameters of θ_S are initialized from the base LLM. Conversely, θ_A are randomly initialized.

Acoustic Tokenizer: We employ XCodec as our acoustic tokenizer which employs a residual vector quantization (RVQ) (Zeghidour et al., 2021) scheme with a codebook size of 1024 across 8 quantization layers. Additionally, we implement the delay pattern. Subsequently, we set the audio boundaries by enclosing this shifted sequence with two special tokens ($\langle | \text{audio_stream_bos} | \rangle$ and $\langle | \text{audio_stream_eos} | \rangle$). Finally, the input representation at each position is obtained by summing the embeddings of the vertically aligned codes across all layers.

5.2 Training Details

In our training pipeline, all models are trained using AdamW with BF16. The details are listed as follows:

Stage 1: We train θ_P and θ_S on InstructS2S-200k for 3 epochs. We use a batch size of 16 and a peak learning rate of 5×10^{-5} . This stage is conducted on 8 NVIDIA H20 GPUs.

Stage 2: In this stage, we freeze θ_S and train θ_A for 2 epochs on Emilia. The batch size is 64, and the learning rate is set to 3×10^{-4} with a warmup of 4,000 steps. This stage is distributed across 32 NVIDIA H20 GPUs.

Stage 3: Finally, we train for 2 epochs with a batch size of 16 and a learning rate of 5×10^{-5} , using 8 NVIDIA H20 GPUs.

5.3 Evaluation on Spoken Question Answering

We evaluate the spoken question answering (SQA) capability of BrainSpeech on three widely used

Table 1: Results on Spoken Question Answering (SQA) benchmarks. We compare BrainSpeech with baselines across different scales. "S→T" denotes Speech-to-Text accuracy, and "S→S" denotes Speech-to-Speech accuracy.

Model	Llama Question		TriviaQA		Web Question	
	S→T	S→S	S→T	S→S	S→T	S→S
MinMo (7B)	78.9	64.1	48.3	37.5	55.0	39.9
GLM-4-Voice (9B)	64.7	50.7	39.1	26.5	32.2	15.9
VITA-Audio (7B)	76.3	64.6	43.6	39.5	44.2	40.0
Moshi (7B)	62.3	21.0	22.8	7.3	26.6	9.2
Llama-Omni (8B)	64.0	45.3	34.2	22.9	31.2	10.7
Llama-Omni2 (7B)	70.3	60.7	38.2	33.5	34.5	31.3
Llama-Omni2 (0.5B)	45.7	38.7	18.5	14.2	17.7	16.8
Llama-Omni2 (1.5B)	62.0	52.7	29.4	24.8	28.2	26.6
Llama-Omni2 (3B)	64.3	55.7	33.6	28.1	30.5	28.0
BrainSpeech (3B)	68.7	63.7	40.5	37.4	39.4	36.0

public benchmarks: Llama Question (Nachmani et al., 2023), TriviaQA (Joshi et al., 2017), and Web Question (Berant et al., 2013).

We benchmark BrainSpeech against a comprehensive suite of representative speechLMs, spanning a wide range of parameter scales (0.5B to 9B). Specifically, we include recent state-of-the-art models such as MinMo (7B) (Chen et al., 2025a), Moshi (7B) (Défossez et al., 2024), GLM-4-Voice (9B) (Zeng et al., 2024) and VITA-Audio (7B) (Long et al., 2025). Furthermore, we compare BrainSpeech with Llama-Omni (8B) (Fang et al., 2024) and Llama-Omni2 family (Fang et al., 2025) (ranging from 0.5B to 7B).

Following established protocols, we employ two evaluation metrics to assess the models:

- *Speech-to-Text (S→T)*: The model receives speech and generates text responses directly. This metric evaluates the pure semantic generation and knowledge retrieval capabilities of the model.
- *Speech-to-Speech (S→S)*: The model receives speech and generates audio responses. The generated speech is then transcribed into text using an ASR model. This metric assesses the capabilities of semantic generation and acoustic quality.

The results are summarized in Table 1. BrainSpeech demonstrates superior performance, achieving the highest accuracy among models of comparable parameter size. Notably, BrainSpeech remains highly competitive even when compared to larger models.

Specifically, these results validate the effectiveness of the architecture of BrainSpeech. By freezing θ_S during stage 2 and 3, BrainSpeech prevents the interference of acoustic gradients, safeguarding the reasoning capabilities of LLMs against catastrophic forgetting, which is the common issue in fine-tuning. Furthermore, the remarkably narrow gap between S→T and S→S metrics attests to the high fidelity of output speech. This confirms that BrainSpeech successfully decouples *what to say* from *how to say*, generating speech with high quality.

5.4 Evaluation on General Capabilities

To further assess the capabilities of BrainSpeech in open-ended scenarios and complex reasoning tasks, we conduct extensive evaluations on subsets of the VoiceBench benchmark (Chen et al., 2024). We confirm the performance across diverse domains:

- *Open-ended QA*: We utilize AlpacaEval and CommonEval to evaluate the ability to handle diverse queries. Following standard protocols, we employ GPT as an impartial judge to assign a score between 1 and 5. For these tasks, we evaluate both audio and text inputs to assess the capabilities of understanding.
- *Reasoning and Knowledge*: We employ MMSU and OpenBookQA to assess multi-choice reasoning capabilities using Accuracy (Acc). For these tasks, we employ Text as input. This design allows us to rigorously evaluate the knowledge retention of θ_S .

Table 2: Evaluation results on VoiceBench subsets. We report Scores (1-5) obtained by DeepSeek (Guo et al., 2025) for AlpacaEval and CommonEval, and Accuracy (%) for MMSU and OpenBookQA. "A" denotes Audio input, and "T" denotes Text input. "-" indicates the modality is unsupported. For MMSU and OpenBookQA, metrics are derived from text input, falling back to the result of audio input only when text input is unsupported.

Model	AlpacaEval		CommonEval		MMSU	OpenBookQA
	A (Score)	T (Score)	A (Score)	T (Score)	Acc	Acc
Moshi (7B)	2.01	-	1.60	-	24.04	26.15
LLaMA-Omni (8B)	3.58	3.70	3.26	3.52	59.01	79.34
VITA-Audio (7B)	2.21	2.73	1.93	2.49	59.41	66.96
GLM-4-Voice (9B)	3.87	-	3.29	-	62.69	44.84
Step-Audio2 (8B)	3.77	3.86	3.14	3.69	66.52	74.89
Mini-Omni (0.5B)	1.95	2.34	2.02	2.55	26.74	30.55
Mini-Omni2 (0.5B)	2.32	2.65	2.18	2.86	27.13	32.09
LLaMA-Omni2 (3B)	3.56	-	3.28	-	21.27	68.06
BrainSpeech (3B)	3.58	3.79	3.15	3.60	53.29	72.31

Baselines: We compare BrainSpeech against a wide array of baselines. We include models such as Mini-Omni/Omni2 (0.5B), LLaMA-Omni2 (3B), and larger-scale models including LLaMA-Omni (8B), Moshi (7B), VITA-Audio (7B), GLM-4-Voice (9B), and Step-Audio 2 (8B).

The quantitative results are presented in Table 2. BrainSpeech exhibits remarkable performance across all tasks. In AlpacaEval and CommonEval, BrainSpeech demonstrates robust instruction-following capabilities. Furthermore, in MMSU and OpenBookQA, BrainSpeech achieves high accuracy, significantly outperforming the comparable LLaMA-Omni2 (3B) and even surpassing larger 7B-scale baselines. This provides strong evidence that our dual-stream architecture, successfully safeguards the model’s semantic generation, preventing the catastrophic knowledge degradation often observed in other speechLMs.

5.5 Evaluation on Speech Synthesis Quality

To evaluate the acoustic quality and expressiveness of BrainSpeech, we conduct assessments on speech generated by audio prompt. The test set aggregates the VoiceBench subsets mentioned above with GenEmotion-en from URO-Bench (Yan et al., 2025), which evaluates the capability to synthesize speech with target emotion based on audio instructions. This ensures a comprehensive evaluation covering both dialogue prosody and specific synthesis capabilities. We employ a combination of objective and subjective metrics:

- *Objective Metrics:* UTMOS (Saeki et al., 2022) is used to assess the generated audio. We also report Content Enjoyment (CE) and Production Quality (PQ) from AudioBox Metrics (Vyas et al., 2023) to assess the overall quality.
- *Subjective Metrics:* Naturalness Mean Opinion Score (NMOS) is used as the metric. We conduct a rigorous human evaluation with 10 listeners. Participants are presented with paired samples (questions and responses) and asked to rate the audio on a scale of 1 to 5 (increments of 0.5 allowed). The criteria strictly focus on prosody, rhythm, and the authenticity of human-like speech (5 = Excellent/Native-like, 1 = Poor/Robotic).

The results are summarized in Table 3. In terms of objective metrics, BrainSpeech achieves highly competitive results, particularly in CE and PQ, surpassing large-scale baselines like GLM-4-Voice. This indicates that θ_A generates high-fidelity speech.

For the subjective evaluation, BrainSpeech outperforms baselines, validating the efficacy of BrainSpeech and the fine-tuning strategy. Our inference strategy effectively resolves boundary artifacts. In contrast, VITA-Audio suffers severely from such discontinuities. Trained on high-fidelity and expressive corpora, BrainSpeech successfully captures the subtle nuances of human articulation, exhibiting superior affective expressiveness and human-level naturalness.

Table 3: Evaluation of speech synthesis quality. We report objective metrics (UTMOS, AudioBox CE/PQ) and subjective human ratings for naturalness (NMOS 1-5). BrainSpeech achieves superior objective and subjective scores.

Model	Objective Metrics			Subjective NMOS
	UTMOS	CE	PQ	
Freeze-Omni (7B)	4.36	6.14	7.75	2.91
LLaMA-Omni (8B)	4.01	6.06	6.95	2.76
VITA-Audio (7B)	3.08	5.76	7.21	1.43
GLM-4-Voice (9B)	4.07	6.20	7.45	3.41
Step-Audio2 (8B)	4.51	6.09	7.66	3.56
Mini-Omni2 (0.5B)	4.44	6.29	7.73	3.08
LLaMA-Omni2 (3B)	4.22	6.34	7.74	3.79
Qwen2.5-Omni (3B)	4.12	6.31	7.71	3.68
BrainSpeech (3B)	4.41	6.31	8.00	3.96

5.6 Ablation Study

To validate the effectiveness of the architecture and training strategy, we examine three variants of models:

- *w/o Acoustic Stream*: A unified model where acoustic attention is removed, retaining only the acoustic FFN and the shared semantic attention. The model is fully fine-tuned on data used in stage 1 to learn understanding and generation simultaneously.
- *Serial Architecture*: Adopts Thinker-Talker where acoustic attention and FFN are appended after the base LLMs. Similar to the above, it is fully fine-tuned on data used in stage 1.
- *BrainSpeech (Stage 2 Only)*: BrainSpeech trained only up to stage 2 (acoustic generation pre-training).

We evaluate these variants on semantic understanding (AlpacaEval/CommonEval Audio Scores) and speech synthesis quality (UTMOS, CE, PQ).

The results are shown in Figure 3. First, regarding architecture, "w/o Acoustic Stream" and "Serial Architecture" suffer significant degradation in semantic performance. This confirms that independent parameter spaces are crucial, and acoustic gradients during full fine-tuning harms the knowledge of base LLMs. Second, regarding the training stage, while stage 2 lays the foundation of audio generation, the full model's performance highlights the necessity of stage 3 for expressive synthesis. It significantly elevates acoustic quality

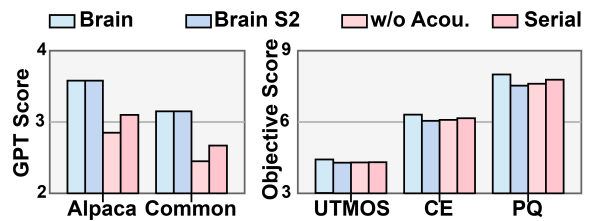


Figure 3: Ablation study results. BrainSpeech best preserves semantics while the fine-tuning stage significantly boosts audio quality.

metrics while maintaining robust semantic alignment established in earlier stages.

6 Conclusion

In this paper, we introduced BrainSpeech, an end-to-end spoken dialogue system that successfully reconciles the conflict between semantic and acoustic generation in SpeechLMs. By employing a biologically inspired dual-stream architecture, BrainSpeech decouples the optimization objectives of *what to say* and *how to say*. This design allows us to safeguard the model's inherent reasoning capabilities and generate high-fidelity speech. Furthermore, our progressive three-stage training paradigm further eliminates the dependency on scarce spoken dialogue data. Extensive evaluations demonstrate that BrainSpeech achieves state-of-the-art performance among models of comparable scale, delivering responses with superior naturalness and emotional expressiveness. We believe BrainSpeech establishes a new baseline for future research into efficient, expressive, and intelligent assistants.

588 Limitations

589 The current implementation prioritizes efficiency
590 with a 3B backbone. While competitive, the po-
591 tential of scaling to larger foundation models re-
592 mains underexplored. We hypothesize that scal-
593 ing up will yield substantial improvements, as
594 larger backbones possess richer semantic repre-
595 sentations. This would facilitate a more compre-
596 hensive capture of the interplay between linguis-
597 tic and paralinguistic features, thereby further ele-
598 vating the quality of both textual and acoustic re-
599 sponses. Future work will investigate these scaling
600 laws to broaden system applicability.

601 References

602 Jonathan Berant, Andrew Chou, Roy Frostig, and Percy
603 Liang. 2013. Semantic parsing on freebase from
604 question-answer pairs. In *Proceedings of the 2013
605 conference on empirical methods in natural lan-
606 guage processing*, pages 1533–1544.

607 Boson AI. 2025. Higgs Audio V2: Redefining Expres-
608 siveness in Audio Generation. [https://github.
609 com/boson-ai/higgs-audio](https://github.com/boson-ai/higgs-audio). GitHub repository.
610 Release blog available at [https://www.boson.ai/
611 blog/higgs-audio-v2](https://www.boson.ai/blog/higgs-audio-v2).

612 Qian Chen, Yafeng Chen, Yanni Chen, Mengzhe Chen,
613 Yingda Chen, Chong Deng, Zhihao Du, Ruize
614 Gao, Changfeng Gao, Zhifu Gao, and 1 others.
615 2025a. Minmo: A multimodal large language
616 model for seamless voice interaction. *arXiv preprint
617 arXiv:2501.06282*.

618 Wenxi Chen, Ziyang Ma, Ruiqi Yan, Yuzhe Liang, Xi-
619 quan Li, Ruiyang Xu, Zhikang Niu, Yanqiao Zhu,
620 Yifan Yang, Zhanxun Liu, and 1 others. 2025b.
621 Slam-omni: Timbre-controllable voice interaction
622 system with single-stage training. In *Findings of
623 the Association for Computational Linguistics: ACL
624 2025*, pages 2262–2282.

625 Yiming Chen, Xianghu Yue, Chen Zhang, Xiaoxue
626 Gao, Robby T. Tan, and Haizhou Li. 2024.
627 Voicebench: Benchmarking llm-based voice assis-
628 tants. *arXiv preprint arXiv:2410.17196*.

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

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

638 Qingkai Fang, Yan Zhou, Shoutao Guo, Shaolei
639 Zhang, and Yang Feng. 2025. Llama-omni2: Llm-
640 based real-time spoken chatbot with autoregres-
641 sive streaming speech synthesis. *arXiv preprint
642 arXiv:2505.02625*.

643 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,
644 Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong
645 Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.
646 Deepseek-r1: Incentivizing reasoning capability in
647 llms via reinforcement learning. *arXiv preprint
648 arXiv:2501.12948*.

649 Haorui He, Zengqiang Shang, Chaoren Wang, Xuyuan
650 Li, Yicheng Gu, Hua Hua, Liwei Liu, Chen
651 Yang, Jiaqi Li, Peiyang Shi, and 1 others. 2024.
652 Emilia: An extensive, multilingual, and diverse
653 speech dataset for large-scale speech generation. In
654 *2024 IEEE Spoken Language Technology Workshop
655 (SLT)*, pages 885–890. IEEE.

656 Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke
657 Zettlemoyer. 2017. Triviaqa: A large scale distantly
658 supervised challenge dataset for reading comprehen-
659 sion. *arXiv preprint arXiv:1705.03551*.

660 KimiTeam, Ding Ding, Zeqian Ju, Yichong Leng,
661 Songxiang Liu, Tong Liu, Zeyu Shang, Kai Shen,
662 Wei Song, and Xu Tan. 2025. Kimi-audio technical
663 report. *arXiv preprint arXiv:2504.18425*.

664 Qingyu Liu, Yushen Chen, Zhikang Niu, Chunhui
665 Wang, Yunting Yang, Bowen Zhang, Jian Zhao,
666 Pengcheng Zhu, Kai Yu, and Xie Chen. 2025.
667 Cross-lingual f5-tts: Towards language-agnostic
668 voice cloning and speech synthesis. *arXiv preprint
669 arXiv:2509.14579*.

670 Zuwei Long, Yunhang Shen, Chaoyou Fu, Heting Gao,
671 Lijiang Li, Peixian Chen, Mengdan Zhang, Hang
672 Shao, Jian Li, Jinlong Peng, and 1 others. 2025.
673 Vita-audio: Fast interleaved cross-modal token gen-
674 eration for efficient large speech-language model.
675 *arXiv preprint arXiv:2505.03739*.

676 Eliya Nachmani, Alon Levkovitch, Roy Hirsch, Ju-
677 lian Salazar, Chulayuth Asawaroengchai, Soroosh
678 Mariooryad, Ehud Rivlin, RJ Skerry-Ryan, and
679 Michelle Tadmor Ramanovich. 2023. Spoken
680 question answering and speech continuation us-
681 ing spectrogram-powered llm. *arXiv preprint
682 arXiv:2305.15255*.

683 Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-
684 man, Christine McLeavey, and Ilya Sutskever. 2023.
685 Robust speech recognition via large-scale weak su-
686 pervision. In *International conference on machine
687 learning*, pages 28492–28518. PMLR.

688 Takaaki Saeki, Detai Xin, Wataru Nakata, Tomoki
689 Koriyama, Shinnosuke Takamichi, and Hiroshi
690 Saruwatari. 2022. Utmos: Utokyo-sarulab sys-
691 tem for voicemos challenge 2022. *arXiv preprint
692 arXiv:2204.02152*.

693	Preeti Saini and Parneet Kau. 2013. Automatic speech recognition: A review. <i>Springer US</i> .	748
694		749
695	Yun Da Tsai, Yu Che Tsai, Bo Wei Huang, Chun Pai Yang, and Shou De Lin. 2023. Automl-gpt: Large language model for automl. <i>arXiv preprint arXiv:2309.01125</i> .	750
696		751
697		752
698		753
699	Apoorv Vyas, Bowen Shi, Matthew Le, Andros Tjandra, Yi-Chiao Wu, Baishan Guo, Jiemin Zhang, Xinyue Zhang, Robert Adkins, William Ngan, and 1 others. 2023. Audiobox: Unified audio generation with natural language prompts. <i>arXiv preprint arXiv:2312.15821</i> .	754
700		755
701		756
702		757
703		758
704		759
705		760
706	Peng Wang, Songshuo Lu, Yaohua Tang, Sijie Yan, Wei Xia, and Yuanjun Xiong. 2024a. A full-duplex speech dialogue scheme based on large language model. <i>Advances in Neural Information Processing Systems</i> , 37:13372–13403.	
707		
708		
709		
710	Xiong Wang, Yangze Li, Chaoyou Fu, Yunhang Shen, Lei Xie, Ke Li, Xing Sun, and Long Ma. 2024b. Freeze-omni: A smart and low latency speech-to-speech dialogue model with frozen llm. <i>arXiv preprint arXiv:2411.00774</i> .	
711		
712		
713		
714		
715	Boyong Wu, Chao Yan, Chen Hu, Cheng Yi, Chengli Feng, Fei Tian, Feiyu Shen, Gang Yu, Haoyang Zhang, and Jingbei Li. 2025. Step-audio 2 technical report. <i>arXiv preprint arXiv:2507.16632</i> .	
716		
717		
718		
719	Zhifei Xie and Changqiao Wu. 2024a. Mini-omni: Language models can hear, talk while thinking in streaming. <i>arXiv preprint arXiv:2408.16725</i> .	
720		
721		
722	Zhifei Xie and Changqiao Wu. 2024b. Mini-omni2: Towards open-source gpt-4o model with vision, speech and duplex. <i>arXiv preprint arXiv:2410.11190</i> .	
723		
724		
725		
726	Ruiqi Yan, Xiquan Li, Wenxi Chen, Zhikang Niu, Chen Yang, Ziyang Ma, Kai Yu, and Xie Chen. 2025. Uro-bench: A comprehensive benchmark for end-to-end spoken dialogue models. <i>arXiv preprint arXiv:2502.17810</i> .	
727		
728		
729		
730		
731	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 22 others. 2024. Qwen2.5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	
732		
733		
734		
735		
736		
737		
738	Guanrou Yang, Chen Yang, Qian Chen, Ziyang Ma, Wenxi Chen, Wen Wang, Tianrui Wang, Yifan Yang, Zhikang Niu, Wenrui Liu, and 1 others. 2025. Emovoice: Llm-based emotional text-to-speech model with freestyle text prompting. <i>arXiv preprint arXiv:2504.12867</i> .	
739		
740		
741		
742		
743		
744	Zhen Ye, Peiwen Sun, Jiahe Lei, Hongzhan Lin, Xu Tan, Zheqi Dai, Qiuqiang Kong, Jianyi Chen, Jiahao Pan, Qifeng Liu, and 1 others. 2025. Codec does matter: Exploring the semantic shortcoming of	
745		
746		
747		
	codec for audio language model. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 25697–25705.	
	Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. 2021. Soundstream: An end-to-end neural audio codec. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 30:495–507.	
	Aohan Zeng, Zhengxiao Du, Mingdao Liu, Kedong Wang, Shengmin Jiang, Lei Zhao, Yuxiao Dong, and Jie Tang. 2024. Glm-4-voice: Towards intelligent and human-like end-to-end spoken chatbot. <i>arXiv preprint arXiv:2412.02612</i> .	

A Formatting Examples of Data

This section illustrates the interleaved dialogue template employed by BrainSpeech. As demonstrated in the listing below, the assistant’s response follows a strict alternating structure of text and audio segments. To ensure fine-grained modality alignment, the text response is segmented into segments of uniform length (text_sequence). Correspondingly, audio segments introduced by the `<|audio_out_bos|>` token contain a **fixed number** of acoustic tokens, maintaining a constant text-to-audio ratio. The special token `<|audio_last_out_bos|>` is reserved for the final audio segment, which encapsulates all **residual acoustic tokens** ensuring the completeness of the speech generation.

Speech QA Interleaved Data Format

```
{
  "messages": [
    {
      "role": "user",
      "content": "<|audio_bos|> <|
        AUDIO|> <|audio_eos|>"
    },
    {
      "role": "assistant",
      "content": "text_sequence_1 <|
        audio_out_bos|>
<|AUDIO_OUT|> <|audio_eos|>
        text_sequence_2
<|audio_out_bos|> <|AUDIO_OUT|> <|
        audio_eos|> text_sequence_3
<|audio_last_out_bos|> <|AUDIO_OUT
        |> <|audio_eos|>"
    }
  ]
}
```

B Dataset Statistics

Dataset	Samples	Avg.Words	Avg.AudioLen
Llama Question	300	8.46	3.00
TriviaQA	1000	12.02	4.72
Web Question	2032	6.76	2.33
AlpacaEval	199	16.32	5.67
CommonEval	200	8.06	4.83
OpenBookQA	455	44.28	18.89
MMSU	3074	53.16	23.61
GenEmotion-en	54	15.35	4.93

Table 4: Statistical information of the datasets used in experiments.

This appendix provides statistical details of the datasets utilized in the experiments. "Samples" in-

dicates the number of samples, "Average Number of Words" (Avg.Words) denotes the average word count of the input of text modality, and "Average Audio Length" (Avg.AudioLen) refers to the average duration of the input of audio modality (unit: seconds).

C Audio Generation with Minimax-Speech-02: Parameter Inference via Qwen2.5

To generate high-quality and expressive audio using Minimax-Speech-02, we leverage its API for audio synthesis. A key step in this process involves inferring critical acoustic and paralinguistic parameters that shape the naturalness of the generated speech. Specifically, we employ Qwen2.5 to analyze input text and predict implicit speech attributes, including pitch, volume, speed, emotion, and gender. Then, these parameters are fed into Minimax-Speech-02 to guide audio generation.

C.1 Parameter Inference Protocol with Qwen2.5

We design a structured instruction for Qwen2.5 to ensure consistent and standardized parameter prediction. The instruction specifies the requirements for each parameter, including valid value ranges, emotion categories, and gender options, as follows:

- **Pitch:** Classified on a scale of $[-3, 3]$, where negative values indicate lower pitch and positive values indicate higher pitch relative to a neutral baseline.
- **Volume:** Rated on a scale of $[0, 4]$, with 0 representing the softest volume and 4 representing the loudest.
- **Speed:** Normalized on a scale of $[0.5, 2.0]$, where 0.5 denotes half the average speech speed and 2.0 denotes double the average speed.
- **Emotion:** Restricted to seven categorical labels: happy, sad, angry, fearful, disgusted, surprised, or calm.
- **Gender:** Specified as male, female or neutral to align with typical speech synthesis gender configurations.

Qwen2.5 processes the input text alongside the instruction and returns the predicted parameters

827 in a machine-readable format, which is then inte-
828 grated into the Minimax-Speech-02 API call.

829 C.2 Example of Parameter Inference and 830 API Input Format

831 The following illustrates the instruction, input text,
832 and corresponding API-compatible format for pa-
833 rameterized audio generation:

```
834 Example of Parameter Inference and  
API Input Format  
{  
  "instruction": "Analyze the following  
text and infer the likely implicit pitch,  
volume, speed, and emotion that would be  
used if the text were spoken aloud.  
Classify pitch on scale [-3, 3], volume on  
scale [0, 4] and speed on scale [0.5, 2.0].  
Determine emotion as one of the  
following: happy, sad, angry, fearful,  
disgusted, surprised, calm. Give the  
gender like male or female. Return the  
result in the format: {'pitch': , 'volume': ,  
'speed': , 'emotion': , 'gender': }",  
  "input": "How can cross training benefit  
groups like runners, swimmers, or  
weightlifters?",  
  "output": ques_0"  
}
```

835 This approach ensures that the generated audio
836 not only accurately conveys the semantic content
837 of the input text but also exhibits natural prosody
838 and emotional alignment, enhancing the overall
839 expressiveness and usability of the synthesized
840 speech.

841 D GPT Evaluation Prompt for 842 VoiceBench

843 This appendix presents the full prompt template
844 used to instruct GPT in automatically evaluat-
845 ing the quality of text responses of models. The
846 prompt standardizes the evaluation criteria for
847 speech interaction scenarios, focusing on rele-
848 vance, accuracy, and conciseness of responses on
849 a 1–5 scale (1 = poorest, 5 = optimal). The prompt
850 (with placeholders {prompt} and {response} re-
851 placed by transcribed user instructions and model
852 responses) is:

GPT Evaluation Prompt

I need your help to evaluate the performance of several models in the speech interaction scenario. The models will receive a speech input from the user, which they need to understand and respond to with a speech output. Your task is to rate the model’s responses based on the provided user input transcription [Instruction] and the model’s output transcription [Response]. Please evaluate the response on a scale of 1 to 5:
1 point: The response is largely irrelevant, incorrect, or fails to address the user’s query. It may be off-topic or provide incorrect information.
2 points: The response is somewhat relevant but lacks accuracy or completeness. It may only partially answer the user’s question or include extraneous information.
3 points: The response is relevant and mostly accurate, but it may lack conciseness or include unnecessary details that don’t contribute to the main point.
4 points: The response is relevant, accurate, and concise, providing a clear answer to the user’s question without unnecessary elaboration.
5 points: The response is exceptionally relevant, accurate, and to the point. It directly addresses the user’s query in a highly effective and efficient manner, providing exactly the information needed.
Below are the transcription of user’s instruction and models’ response:
[Instruction]: prompt
[Response]: response
After evaluating, please output the score only without anything else. You don’t need to provide any explanations.

854 The prompt was applied uniformly across Al-
855 pacEval and CommonEval subsets to ensure con-
856 sistent evaluation of model performance in these
857 scenarios.

E Additional Experiments about Latency

This appendix presents an analysis of latency performance between BrainSpeech and state-of-the-art speechLMs, focusing on two critical metrics for real-time interaction: the latency to generate the first text token (from user input to model’s first text output) and the latency to fully generate the first audio segment (from input to the first speech segment). We compare BrainSpeech against four baselines: GLM-4-Voice (9B), LLaMA-Omni (8B), LLaMA-Omni2 (3B), and VITA-Audio (7B).

Model	FTT (ms)	FAS (ms)
LLaMA-Omni (8B)	289.46	346.73
LLaMA-Omni2 (3B)	204.79	567.84
VITA-Audio (7B)	414.69	1215.89
GLM-4-Voice (9B)	422.77	1562.81
BrainSpeech (3B)	203.75	1050.21

Table 5: Latency comparison of first token/segment generation (unit: milliseconds). "FTT" refers to "Latency of First Text Token", while "FAS" refers to "Latency of First Audio Segment".

As summarized in Table 5, BrainSpeech outperforms GLM-4-Voice (9B) and VITA-Audio (7B) in the two latency metrics. While BrainSpeech exhibits higher latency than LLaMA-Omni (8B) and LLaMA-Omni2 (3B), it achieves significantly superior speech quality compared to these two models (as validated in Table 3). Notably, the higher "FAS" of BrainSpeech relative to LLaMA-Omni2 (3B) stems from the per-segment audio token count: BrainSpeech generates 40 tokens per audio segment, whereas LLaMA-Omni2 (3B) uses only 10 tokens per segment. In practice, we can also reduce the number of audio tokens per segment for BrainSpeech to obtain lower "FAS" latency if needed. Additionally, the latency of BrainSpeech is sufficiently low to support real-time audio interaction already: users only experience a minimal wait for the first audio segment, with no additional delays for subsequent content generation.

F More Discussion About BrainSpeech

In this section, we further elaborate on the broader implications of BrainSpeech beyond its empirical performance presented in the main paper.

As highlighted earlier, BrainSpeech achieves high-quality text responses and more natural,

higher-fidelity audio outputs without relying on scarce multimodal dialogue data. Instead, it leverages only easily accessible resources: pure text dialogue data, synthesized TTS data, and TTS datasets. However, the value of BrainSpeech extends beyond these.

Notably, BrainSpeech aligns more closely with the working mechanism of the brain of human: different brain regions collaborate to perform various tasks, each specializing in its domain, rather than using the whole region (which risks interference between distinct types of information). We can analogize the parameters responsible for text processing in BrainSpeech to the "logical center" of the brain of human: text is a representation constructed by human, rather not naturally occurring in the physical world, and semantic processing based on text enables us to understand the world systematically.

By separating this text (logical) module from the audio module, BrainSpeech mitigates the cross-modal interference that often arises when a single parameter space is forced to handle multiple modalities simultaneously. This modular design also shows scalability to more modalities (e.g., images, video): instead of mixing modal information in a shared parameter space, each modality can be processed by a specialized sub-module that collaborates with the core text logical center. This preserves the clarity of each modality’s processing while enabling coherent multi-modal interaction, extending the strengths of BrainSpeech beyond text-audio dialogue to broader multi-modal scenarios.