# FROM TURN-TAKING TO SYNCHRONOUS DIALOGUE: A SURVEY OF FULL-DUPLEX SPOKEN LANGUAGE MODELS

*Yuxuan Chen*[1*]      *Haoyuan Yu*[2]

[1] Jilin University, Changchun, China [2] Hunan University, Changsha, China
yxchen5522@jlu.edu.cn   y15352176976@hnu.edu.cn

## ABSTRACT

True Full-Duplex (TFD) voice communication—enabling simultaneous listening and speaking with natural turn-taking, overlapping speech, and interruptions—represents a critical milestone toward human-like AI interaction. This survey comprehensively reviews Full-Duplex Spoken Language Models (FD-SLMs) in the LLM era. We establish a taxonomy distinguishing Engineered Synchronization (modular architectures) from Learned Synchronization (end-to-end architectures), and unify fragmented evaluation approaches into a framework encompassing Temporal Dynamics, Behavioral Arbitration, Semantic Coherence, and Acoustic Performance. Through comparative analysis of mainstream FD-SLMs, we identify fundamental challenges—synchronous data scarcity, architectural divergence, and evaluation gaps—providing a roadmap for advancing human-AI communication.

For code and further details, please refer to GitHub[1].

***Index Terms***— True Full-Duplex, Full-Duplex Spoken Language Models, Cognitive Parallelism, Synchronization

## 1. INTRODUCTION

Contemporary SLMs fundamentally lack simultaneous listening and speaking capabilities essential for natural conversation. While LLMs have revolutionized language understanding [1, 2], their spoken dialogue implementations remain constrained by sequential listen-think-speak cycles. Current systems achieve only pseudo-full-duplex (PFD) behavior through time-division multiplexing, failing to match human conversational dynamics [3, 4] characterized by natural turn-taking behaviors illustrated in Fig. 1.

FD-SLMs transform this paradigm from sequential to parallel cognitive architectures. Unlike PFD systems that alternate between listening and speaking, FD-SLMs enable simultaneous encoding and generation within unified processing cycles, supporting natural conversational events including interruptions, backchanneling, and adaptive turn-taking through bidirectional information flow.
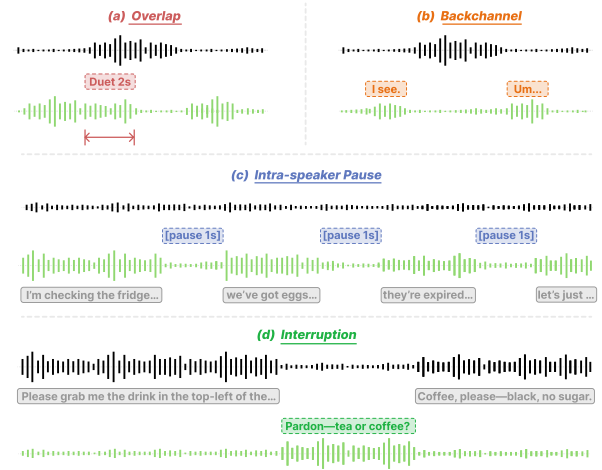
---

[1] https://github.com/elpsykongloo/FD-SLMs



**Fig. 1**. Natural conversations contain turn-taking events: (a) Overlap, (b) Backchannel, (c) Pause, and (d) Interruption.

Early systems demonstrated incremental processing [5] and finite-state control [6], achieving responsiveness without semantic flexibility. LLM integration yielded engineered synchronization [7–9] and end-to-end architectures. Following dGSLM's emergent turn-taking discovery [10], recent advances include hierarchical multi-stream processing [11], next-token-pair prediction (NTPP) [12], and continuous-discrete alignment [13].

Despite these advances, existing surveys [14, 15] treat full-duplex as implementation detail rather than fundamental requirement, lacking systematic FD-SLM design analysis. Evaluation also remains fragmented [16–18].

This paper makes the following primary contributions:

1. **Formal duplex characterization:** Mathematical definitions rigorously distinguish half-duplex, pseudo-full-duplex, and true full-duplex systems, exposing computational requirements for cognitive parallelism.

2. **Architectural taxonomy:** Systematic categorization reveals the design space along synchronization strategy, state management, and training paradigm axes, identifying trade-offs and unexplored opportunities.
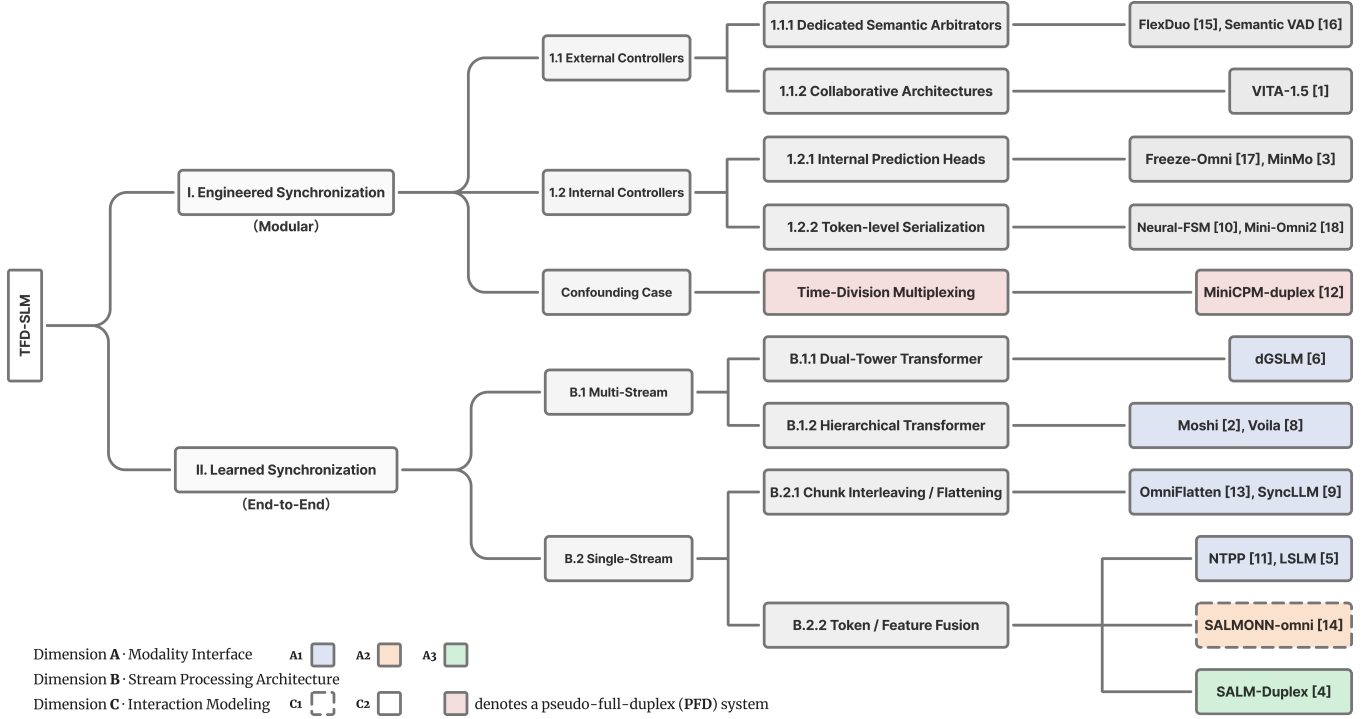
**Fig. 2**. Architectural Taxonomy of Open-source FD-SLMs.

3. **Systematic evaluation and analysis:** We compare streaming architectures across seven FD-SLMs, identify critical data bottlenecks, and establish a four-pillar benchmarking taxonomy revealing the fundamental latency-quality trade-off that constrains current systems.

## 2. FORMALIZATION

FD-SLMs realize cognitive parallelism by concurrently encoding inputs and decoding outputs, enabling real-time output adaptation. Let agent $\mathcal{A}$ interact with environment $\mathcal{E}$. Because direct modeling of continuous audio $X(t)$ is intractable, a discretizer $\mathcal{T}$ yields aligned sequences $S^{\mathcal{E}} = (e_1, \ldots, e_T)$ and $S^{\mathcal{A}} = (a_1, \ldots, a_T)$, aligning $e_t$ with $a_t$ for synchronous interaction.

### 2.1. Joint Probability Perspective

The interaction modeling paradigm models joint distribution $P(S^{\mathcal{E}}, S^{\mathcal{A}})$:

$$P(S^{\mathcal{E}}, S^{\mathcal{A}}) = \prod_{t=1}^{T} P\big(e_t, a_t \mid S^{\mathcal{E}}_{<t}, S^{\mathcal{A}}_{<t}\big). \quad (1)$$

This underlies NTPP [12], simultaneously predicting $(e_t, a_t)$ pairs in decoder-only transformers:

$$\mathcal{L}_{\text{NTPP}}(\theta) = \mathbb{E}_{(S^{\mathcal{E}}, S^{\mathcal{A}})}\Big[\textstyle\sum_{t=1}^{T} \log P\big(e_t, a_t \mid S^{\mathcal{E}}_{<t}, S^{\mathcal{A}}_{<t}; \theta\big)\Big] \quad (2)$$

Earlier approaches [10] approximate through conditional independence with cross-attention, optimizing summed conditional log-likelihoods rather than true joint likelihood.

### 2.2. Conditional Probability Perspective

For interactive agents, the objective becomes modeling $P(S^{\mathcal{A}} \mid S^{\mathcal{E}})$:

$$a_t \sim P\big(a_t \mid S^{\mathcal{E}}_{\leq t}, S^{\mathcal{A}}_{<t}; \theta\big). \quad (3)$$

**Concurrency.** Computing $a_t$ while ingesting $e_{t+1}, \ldots$ requires parallel encoding–decoding [11, 19].

**Real-time constraint.** $\text{Time}(\text{Compute}(a_t)) < 200\,\text{ms}$ [3].

**Self-conditioning.** Dependence on $S^{\mathcal{A}}_{<t}$ ensures coherence and enables echo cancellation [13].

The training objective:

$$\mathcal{L}_{\text{Cond}}(\theta) = \mathbb{E}_{(S^{\mathcal{E}}, S^{\mathcal{A}})}\Big[\textstyle\sum_{t=1}^{T} \log P\big(a_t \mid S^{\mathcal{E}}_{\leq t}, S^{\mathcal{A}}_{<t}; \theta\big)\Big] \quad (4)$$

Training on synchronous data with either objective enables turn-taking dynamics to emerge without supervision.

### 2.3. Hierarchical and Predictive Mechanisms

**Hierarchical Generation.** Systems like Moshi [11] leverage text representations $T^{\mathcal{A}}$:

$$P(S^{\mathcal{A}} \mid S^{\mathcal{E}}) = \int P\big(S^{\mathcal{A}} \mid T^{\mathcal{A}}, S^{\mathcal{E}}\big) \, P\big(T^{\mathcal{A}} \mid S^{\mathcal{E}}\big) \, dT^{\mathcal{A}} \quad (5)$$

A two-stage process generates text tokens ("inner monologue") then audio tokens, merging text-based reasoning with full-duplex capabilities.

**Predictive Synchronization.** SyncLLM [20] predicts upcoming user segments to minimize latency:

$$\hat{e}_{t+1} \sim P\big(\cdot \mid S^{\mathcal{E}}_{\leq t}, S^{\mathcal{A}}_{\leq t}\big), \quad a_{t+1} \sim P\big(\cdot \mid S^{\mathcal{E}}_{\leq t}, \hat{e}_{t+1}, S^{\mathcal{A}}_{\leq t}\big) \quad (6)$$

**Table 1**. Comparative analysis of architectural components in open-source FD-SLMs.

| Model | Input Perception | Core Processing | Output Synthesis |
|---|---|---|---|
| dGSLM | HuBERT + k-means clustering | Two-tower Transformer with cross-attention | HiFi-GAN unit vocoder |
| Moshi | Mimi neural codec (RVQ) | RQ-Transformer joint autoregression | Mimi decoder |
| SyncLLM | HuBERT features | Interleaved and predictive synchronization | HiFi-GAN vocoder |
| SALMONN-omni | Mamba streaming encoder | Dynamic control tokens for stream management | CosyVoice2 with fixed-length generation |
| MinMo | SenseVoice-Large + projector | Full-Duplex Predictor (FDP) head | CosyVoice2 chunk-aware flow-matching |
| FlexDuo | Qwen2-Audio encoder | Finite-state machine control | External TTS delegation |
| VITA 1.5 | Conv + Transformer encoder | Dual LLM instances with shared KV cache | TiCodec decoder |

## 3. TAXONOMY

Cognitive parallelism, enabling simultaneous speech encoding and output decoding, requires departing from sequential Transformer architectures. Figure 2 shows current approaches following two paradigms: **engineered synchronization** via modular architectures and **learned synchronization** through end-to-end systems.

### 3.1. Engineered Synchronization

Modular approaches enhance dialogue engines with specialized components, eliminating retraining through explicit state arbitration. The duplex controller—a neural FSM—extends beyond acoustic VAD to perform semantic arbitration, differentiating interruptions from backchannels and noise.

**External controllers.** External controllers maintain independence from the core engine. FlexDuo introduces a ternary FSM with an idle state for selective attention [7]. Semantic VAD uses lightweight ($\sim$0.5B) models analyzing ASR outputs to minimize computational load [21]. VITA-1.5 employs dual instances that swap roles upon interruption detection, trading computational cost for latency [22].

**Internal controllers.** Internal controllers embed control logic within the engine architecture. Freeze-Omni performs chunk-wise state prediction on frozen LLMs [23]; MinMo's Full Duplex Predictor reads embeddings for turn-yielding decisions [24]. Neural-FSM extends vocabularies with FSM tokens enabling autonomous state management through next-token prediction [8]. Mini-Omni2 implements command-based interruption via semantic state tokens [24].

### 3.2. Learned Synchronization

End-to-end architectures natively process bidirectional audio streams. Following dGSLM's demonstration of emergent turn-taking from raw audio [10], these systems make full-duplex capabilities intrinsic. The challenge lies in reconciling Transformer sequentiality with conversational parallelism.

**Modal interfaces.** Modal interfaces vary in representation. Codec-based approaches [10–12, 20, 25] discretize audio into tokens despite sequence elongation. SALMONN-omni directly processes continuous embeddings [13]. SALM-Duplex combines continuous inputs with discrete outputs for an accuracy–latency tradeoff [26].

**Stream processing.** Stream processing follows multi-stream or single-stream paradigms. Multi-stream approaches like dual-tower architectures use cross-attention for synchronization [10], while Moshi's RQ-Transformer jointly models user/agent audio and internal monologue [11]. Single-stream methods serialize inputs for standard decoders: SyncLLM interleaves chunks with synchronization tokens [20], NTPP uses pairwise causal masking [12], and LSLM/SALM-Duplex explore varying fusion depths [19, 26].

**Interaction modeling.** Interaction modeling predominantly employs implicit dynamics where models control turn-taking through silence or audible token generation without explicit supervision [10–12, 20, 25]. In contrast, SALMONN-omni's Dynamic Thinking mechanism [13] generates control tokens for explicit state management, positioning the LLM as the duplex predictor within an end-to-end framework.

## 4. EVALUATION

FD-SLMs demand coordinated assessment across three interdependent axes: streaming architectures enabling real-time interaction, conversational training data, and comprehensive benchmarking methodologies.

### 4.1. Architectural Components

FD-SLMs require specialized streaming architectures achieving sub-200 ms latency for natural turn-taking [14, 15]. Table 1 summarizes strategies across three critical stages.

**Input Perception.** Continuous encoding with minimal lookahead is essential. While conventional encoders need causal adaptation, purpose-built streaming encoders operate natively [13, 27, 28]. Discrete paradigms employ strictly causal/near-zero-lookahead neural codecs [11, 12]; tokenizer chunk granularity fundamentally bounds perception latency [12, 29–31].

**Core Processing.** Concurrent streams are synchronized via cross-attention [10], joint autoregression [11], predictive synchronization [20], or explicit control mechanisms [7, 24]. A

**Table 2**. Comprehensive evaluation of representative open-source FD-SLMs across four dimensions.

| Model | Temporal Dynamics | | Behavioral Arbitration | | | Semantic Coherence | | Acoustic Performance | |
|---|---|---|---|---|---|---|---|---|---|
| | FTO ($\downarrow$) | SL ($\downarrow$) | IRD ($\downarrow$) | ISR ($\uparrow$) | WER ($\downarrow$) | PPL ($\downarrow$) | QA Acc ($\uparrow$) | N-MOS ($\uparrow$) | M-MOS ($\uparrow$) |
| Human | $\sim$0.20 s | $\sim$0.30 s | 2.32 s | 93.69% | 1.5% | 10.2 | 92% | 4.92 ($\pm$0.02) | 4.85 ($\pm$0.03) |
| dGSLM | 0.33 s ($\pm$0.12) | 0.15 s ($\pm$0.03) | 1.33 s | 60.31% | 25% ($\pm$3.4) | 334.4 | 17.2% | 3.85 ($\pm$0.12) | 1.38 ($\pm$0.10) |
| NTPP | 0.30 s ($\pm$0.15) | 0.18 s ($\pm$0.05) | 1.30 s | 80.82% | 7.5% ($\pm$1.22) | 35 | 55.2% | 4.15 ($\pm$0.06) | 3.95 ($\pm$0.04) |
| Moshi | 2.22 s ($\pm$0.70) | 0.75 s ($\pm$0.10) | 1.44 s | 77.73% | 5.20% ($\pm$0.13) | 59.3 | 33.8% | 3.90 ($\pm$0.07) | 3.75 ($\pm$0.06) |
| SALMONN-omni | 0.38 s ($\pm$0.10) | 0.25 s ($\pm$0.08) | 1.38 s | 85.6% | 8.40% ($\pm$0.20) | 21.1 | 61% | 3.85 ($\pm$0.10) | 3.95 ($\pm$0.15) |
| VITA-1.5 | 2.10 s ($\pm$0.65) | 0.12 s ($\pm$0.05) | 9.49 s | 78.53% | 5.45% ($\pm$0.10) | 26.8 | 50.5% | 4.00 ($\pm$0.08) | 4.10 ($\pm$0.10) |
| Freeze-Omni | $-$0.40 s ($\pm$0.05) | 1.11 s ($\pm$0.17) | 9.25 s | 54.97% | 7.30% ($\pm$0.05) | 30.2 | 56.9% | 3.80 ($\pm$0.10) | 3.90 ($\pm$0.07) |

100–200 ms "cognitive clock" sets perception–reaction granularity [9, 11, 20]. KV-cache efficiency directly affects sustained responsiveness [12, 32].

**Output Synthesis.** Discrete models reuse codec decoders for minimal latency [33, 34]. Continuous pipelines employ chunk-aware flow-matching [35], fixed-length interleaved generation [13], or tightly coupled LLM–vocoder stacks [22].

### 4.2. Training Data

Data scarcity remains critical: FD-SLMs require synchronized multi-channel spontaneous dialogue, unavailable in monologue corpora. Current training uses limited datasets [10–12] constraining diversity (see Table 3 for examples; full listings in our repository).

**Table 3**. Publicly Available Datasets for FD-SLM

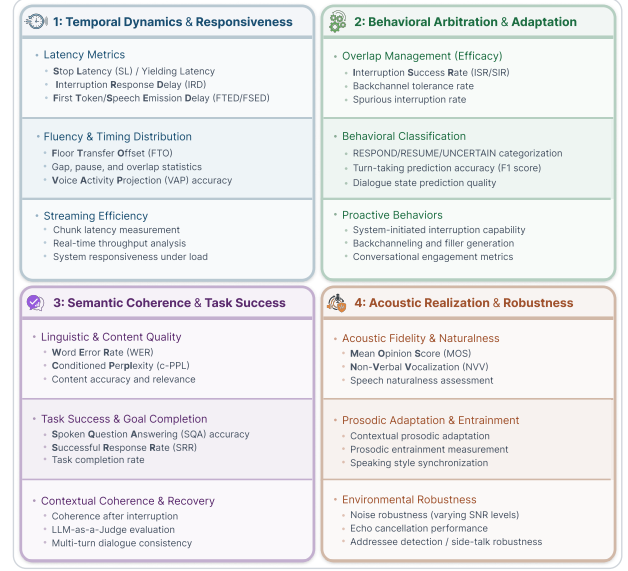| Dataset | Lang | Scene | Channels | Hours |
|---|---|---|---|---|
| AMI Meeting Corpus | EN | meeting | 8 | 100 |
| ICSI Meeting Corpus | EN | meeting | 6 | 70 |
| LibriCSS | EN | meeting | 7 | 10 |
| Fisher English | EN | phone | 2 | 1,960 |
| SEAME (Mandarin–English CS) | EN+ZH | interview | 2 | 192 |
| HKUST Mandarin Telephone | ZH | phone | 2 | 149 |

Synthetic TTS generation [20] fails to capture prosodic entrainment and overlap dynamics, limiting generalization. Progress requires end-to-end conversational synthesis and advanced source separation for single-channel data.

### 4.3. Benchmarking Framework

Conventional metrics built for half-duplex systems [1, 36] fail to capture real-time FD behaviors: when models speak, how they intervene, and conversational floor arbitration [6].

Historical fragmentation through model-specific metrics [8, 10–12] prevented systematic comparison. Recent standardization efforts [16–18] enable reproducible evaluation via our four-pillar taxonomy (Fig. 3).

Table 2 reveals critical gaps: while acoustic quality approaches human levels, temporal dynamics vary widely, behavioral arbitration underperforms (ISR: 54-86% vs. 94% human), and semantic coherence trades off against responsiveness—demonstrating that human-parity FD requires paradigmatic architectural advances.



**Fig. 3**. Four-Pillar Taxonomy of Benchmarking FD-SLMs.

## 5. CONCLUSION

FD-SLMs mark a paradigm shift from turn-based to synchronous dialogue. Through cognitive concurrency formalization and our taxonomy distinguishing Engineered from Learned Synchronization, we clarify fundamental design trade-offs. Our four-pillar evaluation reveals that while acoustic quality approaches human levels, critical gaps persist: inconsistent temporal dynamics, suboptimal behavioral arbitration, and inverse latency-coherence correlation.

Progress requires addressing interconnected challenges. Architectural fragmentation prevents scalable designs aligned with LLM scaling laws. Data scarcity—particularly synchronized multi-channel recordings and non-English resources [37]—constrains learning. Current evaluation lacks proactive behavior metrics [25], while ultra-low latency introduces safety risks requiring real-time filtering.

Advancing FD-SLMs demands architectural convergence, synthetic data capturing authentic dynamics, comprehensive behavioral evaluation, and robust safety mechanisms. Only through coordinated efforts can we achieve truly human-like conversational AI that is responsive, scalable, and ethically deployable.

# 6. REFERENCES

[1] OpenAI et al., "GPT-4o system card," 2024.

[2] Boyong Wu et al., "Step-audio 2 technical report," 2025.

[3] Tanya Stivers et al., "Universals and cultural variation in turn-taking in conversation," *PNAS*, 2009.

[4] Stephen C. Levinson and Francisco Torreira, "Timing in turn-taking and its implications for processing models of language," *Frontiers in Psychology*, 2015.

[5] Antoine Raux and Maxine Eskenazi, "A finite-state turn-taking model for spoken dialog systems," in *Proc. NAACL-HLT*, 2009.

[6] Gabriel Skantze, "Turn-taking in conversational systems and human-robot interaction: A review," *Computer Speech & Language*, 2021.

[7] Ziyang Liao et al., "Flexduo: A pluggable system for enhancing spoken dialogue models with full-duplex capabilities," 2025.

[8] Zheng Wang et al., "Neural-FSM: A full-duplex speech dialogue scheme based on large language model," in *Proc. NeurIPS*, 2024.

[9] Jun Zhang et al., "Omniflatten: A unified framework for spoken language model via progressive flattening," in *Proc. ACL*, 2025.

[10] Anh-Duy Nguyen et al., "Generative spoken dialogue language modeling," *TACL*, 2023.

[11] Alexandre Défossez et al., "Moshi: a speech-text foundation model for real-time dialogue," 2024.

[12] Qichao Wang et al., "NTPP: Generative speech language modeling for dual-channel spoken dialogue via next-token-pair prediction," 2025.

[13] Dong Yu et al., "Salmonn-omni: A codec-free LLM for full-duplex speech understanding and generation," 2025.

[14] Siddhant Arora et al., "On the landscape of spoken language models: A comprehensive survey," 2025.

[15] Jia Cui et al., "Recent advances in speech language models: A survey," 2025.

[16] Yizhou Peng et al., "Fd-bench: A full-duplex benchmarking pipeline designed for full duplex spoken dialogue systems," 2025.

[17] Guan-Ting Lin et al., "Full-duplex-bench: A benchmark to evaluate full-duplex spoken dialogue models on turn-taking capabilities," 2025.

[18] Guan-Ting Lin et al., "Full-duplex-bench v1.5: Evaluating overlap handling for full-duplex speech models," 2025.

[19] Ziyang Ma et al., "Language model can listen while speaking," 2024.

[20] Aditya Veluri et al., "Beyond turn-based interfaces: Synchronous LLMs as full-duplex dialogue agents," 2024.

[21] Mohan Shi et al., "Semantic VAD: Low-latency voice activity detection for speech interaction," 2023.

[22] Chaoyou Fu et al., "Vita-1.5: Towards GPT-4o level real-time vision and speech interaction," 2025.

[23] Xiong Wang et al., "Freeze-omni: A smart and low latency speech-to-speech dialogue model with frozen LLM," 2024.

[24] Zhifei Xie et al., "Mini-omni2: Towards open-source GPT-4o with vision, speech and duplex capabilities," 2024.

[25] Yemin Shi et al., "Voila: A sophisticated, synchronous, and swift spoken language model," 2025.

[26] Ke Hu et al., "Salm-duplex: Efficient and direct duplex modeling for speech-to-speech language model," 2025.

[27] Yangyang Shi et al., "Emformer: Efficient memory transformer based acoustic model for low latency streaming speech recognition," in *Proc. ICASSP*, 2021.

[28] Zengwei Yao et al., "Zipformer: A faster and better encoder for automatic speech recognition," in *Proc. ICLR*, 2024.

[29] Pooneh Mousavi et al., "Discrete audio tokens: More than a survey!," 2025.

[30] Rithesh Kumar et al., "High-fidelity audio compression with improved RVQGAN," in *Proc. NeurIPS*, 2023.

[31] Yuanzhe Xu et al., "Wavtokenizer: A novel general-purpose audio-to-token converter," in *Proc. ICLR*, 2025.

[32] Zhen Li et al., "Connector-s: A survey of connectors in multimodal large language models," 2025.

[33] Zhengyan Sheng et al., "Syncspeech: Low-latency and efficient dual-stream text-to-speech based on temporal masked transformer," 2025.

[34] Sambal Shikhar et al., "LLMVoX: A zero-shot, personalized, and streaming speech synthesis leveraging large language models," in *Findings of ACL*, 2025.

[35] Qian Chen et al., "Minmo: A multimodal large language model for seamless voice interaction," 2025.

[36] Takaaki Saeki et al., "UTMOS: UTokyo-sarulab MOS prediction system for voice conversion challenge 2022," in *Proc. SSW*, 2022.

[37] Shintaro Ohashi et al., "Towards a japanese full-duplex spoken dialogue system: Data collection, modeling, and evaluation," 2025.

[38] OpenBMB et al., "Minicpm-llama3-v 2.5: An 8b-scale multimodal LLM," 2024.

[39] Jing Peng et al., "A survey on speech large language models for understanding," 2024.

[40] Jinglin Chen et al., "Wavchat: A survey of spoken dialogue models," 2024.

[41] Guillermo Castillo-López et al., "A survey of recent advances on turn-taking modeling in spoken dialogue systems," in *Proc. IWSDS*, 2025.

[42] Alexandre Défossez et al., "High-fidelity neural audio compression," *Trans. Mach. Learn. Res.*, 2023.

[43] Hubert Siuzdak, "Vocos: Closing the gap between time-domain and fourier-based neural vocoders for high-quality audio synthesis," in *Proc. ICLR*, 2024.