# $\tau$-voice: Benchmarking Full-Duplex Voice Agents on Real-World Domains

**Anonymous Authors**[1]

## Abstract

Full-duplex voice agents—systems that listen and speak simultaneously—are rapidly moving from research to production. However, existing evaluations address conversational dynamics and task completion in isolation. We introduce $\tau$-voice, a benchmark for evaluating voice agents on grounded tasks with real-world complexity: agents must navigate complex multi-turn conversations, adhere to domain policies, and interact with the environment. The framework extends $\tau^2$-bench into the first voice agent benchmark combining verifiable completion of complex grounded tasks, full-duplex interaction, and realistic audio—enabling direct comparison between voice and text performance. A controllable and realistic voice user simulator provides diverse accents, realistic audio environments, and rich turn-taking dynamics; by decoupling simulation from wall-clock time, the user simulator can use the most capable LLM without real-time constraints. We evaluate task completion (pass@1) and voice interaction quality across 278 tasks: while GPT-5 (reasoning) achieves **80%**, voice agents reach only **29–42%** under clean conditions and **19–30%** under realistic conditions with noise and diverse accents—a **50–61pp gap**; qualitative analysis confirms 75–90% of failures stem from agent behavior, suggesting that observed failures primarily reflect agent behavior under our evaluation setup. $\tau$-voice provides a reproducible testbed for measuring progress toward voice agents that are natural, conversational, and reliable.

## 1. Introduction

The next frontier in conversational AI is **full-duplex voice interaction**—natural spoken conversations where systems listen and speak simultaneously, handle interruptions gracefully, and make real-time turn-taking decisions (Gartner, 2024; 2025; Moore, 2025). Unlike turn-based interactions where users speak, wait, and speak again, full-duplex systems operate in continuous time without explicit turn boundaries.

A new generation of **audio-native language models** enables this vision, processing speech end-to-end without intermediate transcription. Customer service is a primary application: voice remains the preferred channel for complex issues where customers need to explain nuanced problems or resolve urgent matters.

Existing work evaluates whether these models can hold a conversation—but can they *simultaneously* process a return, modify an order, or resolve a billing dispute, with the reliability we expect from text-based agents?

### 1.1. Why End-to-End Evaluation Matters

Voice agents must excel at two capabilities: **task completion** (reasoning about requests, calling tools correctly, modifying database state) and **conversation management** (turn-taking, interruptions, backchanneling in continuous time). Existing benchmarks evaluate each in isolation: $\tau$-bench and $\tau^2$-bench (Yao et al., 2024; Barres et al., 2025) measure tool use on realistic customer service tasks but in text-only, turn-based settings; Full-Duplex-Bench and its v2 (Lin et al., 2025b;a) evaluate turn-taking and interruptions but on synthetic tasks without real tool calls (§2). What remains unexplored is evaluating both together: voice interaction grounded in consequential tasks.

Voice compounds task difficulty in ways text does not. Speech lacks punctuation, contains fillers and disfluencies, and requires verbally encoding special characters. The *audio environment* (background noise, accents, telephony compression) introduces errors that propagate across turns. Real-time *conversational dynamics* (interruptions, backchannels, turn-taking) demand that agents respond fluidly without long silences.

Consider:

> *A customer calls to make changes to their account. Due to background noise and an unfamiliar ac-*

[1]Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

*cent, the agent mishears their name and authentication fails. Does the agent ask them to spell it? If the customer spells it out, does the agent transcribe it correctly despite the noise? If so, does it fix the authentication tool call—or does it make a mistake in combining the information spread across the turns?*

Such failures cannot be captured by evaluating ASR, dialogue state tracking, and tool use separately. They also pose **accessibility concerns**: users with non-standard accents, speech impediments, or noisy environments may be systematically underserved by voice agents that perform well only under ideal conditions.

### 1.2. Our Contributions

We present $\tau$-**voice**, extending $\tau^2$-bench to full-duplex voice interaction:

1. **First voice agent benchmark combining verifiable completion of complex grounded tasks, full-duplex interaction, and realistic audio.** Existing benchmarks evaluate these dimensions in isolation (§2). $\tau$-voice is the first to combine all three and enables direct comparison between voice and text agent performance on grounded tasks.

2. **Controllable and realistic voice user simulator.** A voice user simulator with diverse accents, realistic audio environments, and rich turn-taking dynamics. By decoupling simulation time from wall-clock time, our user simulator can use the most capable LLM without real-time constraints, ensuring reliable instruction following and turn-taking decisions.

3. **Empirical findings.** We benchmark Google, OpenAI, and xAI, ablating acoustic factors (noise, accents, user behaviors). Figure 1 summarizes our headline result:

   • *A large voice-text gap remains*: Even under Clean conditions (clean audio, no interruptions), voice agents achieve only 29–42% vs 80% GPT-5 (reasoning)—a 38–51pp gap.
   • *Realistic audio exacerbates the gap*: Under Realistic conditions (noise, accents, user behaviors), performance falls further to 19–30%. Among factors, accents hurt most ($-13$pp), with potential accessibility implications.
   • *Provider trade-offs*: Google handles individual factors well ($-5$pp combined) but compounds under full Realistic conditions ($-11$pp). For turn-taking, xAI achieves best latency (0.99s) and responsiveness (85%) but interrupts users on average once per turn; OpenAI has best selectivity (74%) but worst latency

(2.22s). No provider masters both task completion and conversational dynamics.
• *Failures are primarily agent errors*: Qualitative analysis of 40 failed tasks confirms that 75–90% of failures stem from agent behavior, suggesting that observed failures primarily reflect agent behavior under our evaluation setup.
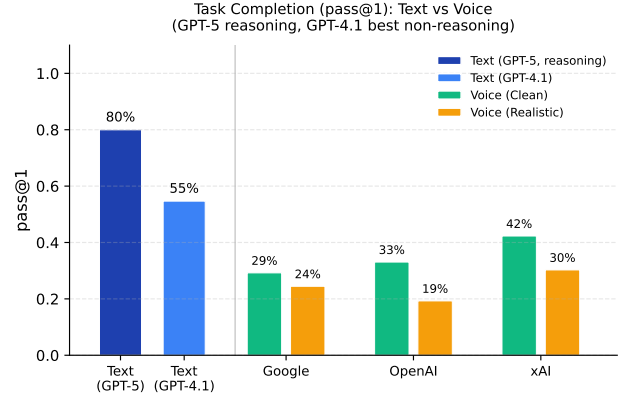


*Figure 1.* Task completion (pass@1) averaged across all domains. GPT-5 (reasoning) achieves 80%. Voice agents show two levels of degradation: under **Clean** conditions (clean audio, no interruptions), performance drops to 29–42% ($-38$ to $-51$pp); under **Realistic** conditions (realistic audio, interruptions), it falls further to 19–30% ($-50$ to $-61$pp from text).

## 2. Related Work

Evaluating voice agents requires measuring both *what* they accomplish and *how* they converse. Table 1 summarizes how existing benchmarks address three key dimensions: **Task Completion** (tasks requiring correct API calls with verifiable database state changes), **Full-Duplex** (simultaneous bidirectional speech with turn-taking and interruptions), and **Realistic Audio Environment** (diverse speaker characteristics, accents, background noise, channel degradation, and disfluencies).

*Table 1.* Comparison of evaluation dimensions across benchmarks. Prior work advances individual dimensions; $\tau$-voice combines all three.

| | Task Completion | Full-Duplex | Realistic Audio Env. |
|---|---|---|---|
| *Task-Oriented (Text)* | | | |
| $\tau$-bench (Yao et al., 2024) | ✓ | | |
| $\tau^2$-bench (Barres et al., 2025) | ✓ | | |
| *Conversational Dynamics* | | | |
| Full-Duplex-Bench (Lin et al., 2025b;a) | ∼ | ✓ | |
| Talking Turns (Arora et al., 2025) | | ✓ | |
| *Speech Understanding* | | | |
| VoiceBench (Chen et al., 2024) | | | ✓ |
| VocalBench (Liu et al., 2026) | | | ✓ |
| Audio MultiChallenge (Gosai et al., 2025) | | | ✓ |
| $\tau$-**voice** | ✓ | ✓ | ✓ |

## 2.1. Task-Oriented Agents (Text)

$\tau$-bench (Yao et al., 2024) evaluates agents on customer service tasks with verifiable database outcomes (§1). $\tau^2$-bench (Barres et al., 2025) extends this to dual-control settings where users also have tool access. Both operate entirely in text—no acoustic variation or real-time constraints.

## 2.2. Conversational Dynamics

Full-Duplex-Bench (Lin et al., 2025b) introduced automatic metrics for pause handling, backchanneling, turn-taking, and interruption management. V2 (Lin et al., 2025a) extends this to multi-turn evaluation with task families (daily scenarios, correction handling, entity tracking, safety) and an automated examiner that enforces staged goals. However, these tasks remain scripted scenarios rather than real tool calls against databases. V2's real-time streaming approach also limits fine-grained control—interruption, backchannel, and yield timing are not precisely configurable. In contrast, our tick-based orchestrator enables configurable turn-taking behavior, making it easy to increase or decrease realism and difficulty. Talking Turns (Arora et al., 2025) evaluates turn-taking using a model trained on human judgments, revealing that current models interrupt inappropriately and rarely backchannel.

## 2.3. Speech & Audio Understanding

VoiceBench (Chen et al., 2024) evaluates ASR robustness across diverse speaker characteristics and acoustic environments. VocalBench (Liu et al., 2026) evaluates vocal conversational abilities—response quality, acoustic performance, and conversational flow. Audio MultiChallenge (Gosai et al., 2025) provides multi-turn context but evaluates only a single model response, testing memory and coherence with disfluencies. Related work addresses prosody, disfluencies, and speaker diversity in natural speech (Zhang et al., 2025; Wang et al., 2025). Beyond robustness, paralinguistic benchmarks (Jiang et al., 2025; Yang et al., 2025; Ao et al., 2025) evaluate understanding of emotion, accent, and prosody. While these benchmarks reveal important capability gaps, they evaluate speech processing in isolation from task completion.

## 2.4. The Missing Intersection

As Table 1 shows, no existing benchmark combines all three dimensions. $\tau$-voice addresses this gap.

## 3. Methods

We extend $\tau^2$-bench to voice interactions through three components: a full-duplex orchestrator enabling reproducible and controllable evaluation, a realistic voice user simulator,

and metrics capturing both task completion and interaction quality.
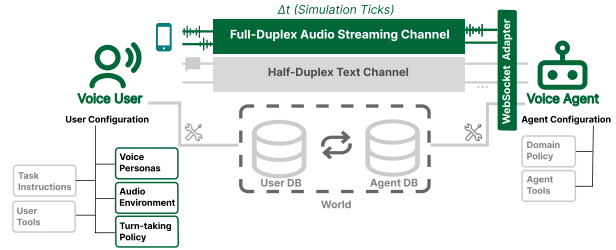
### 3.1. Full-Duplex Orchestrator



*Figure 2.* $\tau$-voice extends $\tau^2$-bench (gray) with voice-specific components (green): a voice user simulator with configurable personas, audio environment, and turn-taking policy; a full-duplex audio streaming channel discretized into simulation ticks; and a provider adapter for adding new voice APIs. Task infrastructure (instructions, tools, databases, domain policies) is inherited.

The orchestrator coordinates the interaction loop between the voice user simulator and the agent API, managing audio exchange, turn-taking events, and evaluation logging. Voice agent APIs (OpenAI Realtime (OpenAI, 2025), Gemini Live (Vertex AI, 2025), xAI Grok (xAI, 2025)) are designed for continuous real-time streaming with bidirectional audio flow and voice activity detection (VAD) for turn-taking. Crucially, these APIs index events on *audio time* rather than wall-clock time—audio can be sent faster or slower than real-time and the API processes it according to audio timestamps.

This decoupling enables our tick-based orchestrator: by advancing simulation time independently of wall-clock time, we allow the user simulator to use the most capable LLM without real-time constraints, ensuring reliable instruction following and turn-taking decisions. This enables reproducibility and fine-grained control over the timing of all turn-taking actions.

**Discrete Simulation Time.** We discretize the continuous audio stream into fixed-duration **ticks** ($\tau = 200ms$ by default). Each tick, both parties exchange exactly $\tau$ ms of audio, enabling true full-duplex interaction where both can speak simultaneously. Since audio generation may not align with tick boundaries, both sides buffer; on interruption, the buffer is cleared, truncating the agent's in-progress response (formal details in Appendix B.1). The agent returns both audio and transcript text each tick, with text distributed proportionally to audio duration (Appendix B.2); overlapping speech is linearized to sequential text for the user simulator LLM (Appendix B.3).

**Controllability and Reproducibility.** Decoupling from real-time enables fine-grained control over all simulation

parameters. Conversational dynamics are configurable: silence thresholds before responding, interruption check intervals, yield timing after overlap. The audio environment is fully parameterized: background noise SNR and drift, burst noise rate and intensity, telephony compression settings, and frame drop probability via a Gilbert-Elliott model. Voice personas specify accent, speaking style, and prosody. This enables systematic ablations isolating the impact of individual factors on task performance. Given a seed, all stochastic elements are deterministic for controlled comparison across agents; full reproducibility is limited only by LLM output variance.

### 3.2. Voice User Simulator

Voice interactions introduce challenges absent from text: the *audio environment* degrades signals, and *conversational dynamics* require real-time turn-taking decisions. Our simulator addresses these by generating realistic caller audio through a pipeline (Figure 3) combining text generation, speech synthesis, audio environment simulation, and conversational dynamics.

To isolate agent performance from transcription artifacts, the simulator receives the agent's transcript directly rather than transcribing agent speech.
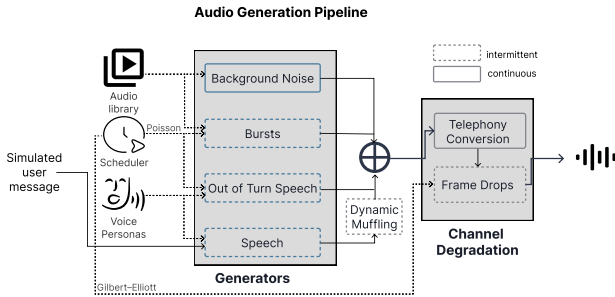


*Figure 3.* Voice user simulator pipeline. Each tick, the simulator generates text, synthesizes speech with a persona, mixes in environmental audio, and applies telephony degradation to produce realistic caller audio.

**Speech Generation.** User simulator prompts produce natural spoken language: disfluencies and fillers ("um", "uh"), verbalized special characters ("at" not "@"), and terse responses. Generated text is synthesized using voice personas—each with a dedicated TTS voice and system prompt guiding speech style and prosody. We define seven personas spanning diverse accents and demographics (Appendix C).

**Audio Environment.** We simulate realistic telephony conditions by mixing synthesized speech with environmental audio: continuous background noise (chatter, traffic) and intermittent bursts (phone rings, dog barks) drawn from

recorded samples. Out-of-turn speech—synthesized phrases like "hold on" and vocal tics like coughs and sneezes—simulates moments when callers are distracted. Effects degrade the signal: dynamic muffling simulates movement away from the microphone, telephony conversion applies G.711 $\mu$-law compression at 8kHz, and frame drops simulate packet loss. All streams are mixed to target signal-to-noise ratios relative to the primary speech. Parameters appear in Appendix D.

**Turn-Taking Policy.** The simulator combines configurable threshold-based timing with LLM-driven decisions. For example, the user waits for a silence threshold (default 1s) before responding. During agent speech, an LLM periodically evaluates whether to interrupt based on conversation context. A separate LLM decides whether to backchannel ("mm-hmm"), and if the agent interrupts, the user yields after a configurable overlap duration. Full prompts appear in Appendix F; Table 2 illustrates these dynamics.

### 3.3. Evaluation

Voice evaluation requires capturing both task outcomes and conversational behavior. We instrument each simulation to log turn-taking events, audio effects, and agent responses, then derive metrics for task success and voice interaction quality.

*Table 2.* Key moments from the Task 41 trajectory (Figure 4). At 8s, the agent interrupts; at 68s the user interrupts and the agent yields but fails to respond for 5 seconds; at 82s the agent incorrectly responds to non-agent-directed speech [in brackets]; at 113s the user interrupts but the agent does not yield; at 121s the agent correctly continues through a backchannel.

| Time | User | Agent | Event |
|---|---|---|---|
| 5–8s | Hi, I have two prob- | | |
| 8s | -lems. First, I ordered | Hello! | *agent int.* |
| | | … | |
| 60–67s | | ...Which would you like | |
| 67–68s | | to do first? | |
| 68–69s | Jigsaw first. | | *user int., yield* |
| 69–74s | | | *no response* |
| 74–77s | Can you switch it... | | |
| | | … | |
| 77–82s | | To confirm, you want to | |
| 82s | [Give me a moment.] | exchange the puzzle– | *non-dir., yield* |
| 84s | | Sure, take your time. | *error: responds* |
| | | … | |
| 108–113s | | ...on order #W4082615. Is | |
| 113–114s | Yeah, that's it. | that the one? We can exch- | *user int.* |
| 114–115s | | -ange it for a puzzle... | *no yield* |
| | | … | |
| 115–121s | | ...500-piece puzzles. Wo- | |
| 121–122s | mm-hmm | -uld you like to exchange | *backchannel* |
| 122–128s | | it for one of those? | *continues* |

**Timeline Walkthrough.** Figure 4 illustrates our evaluation on a 3-minute Retail conversation with street noise.
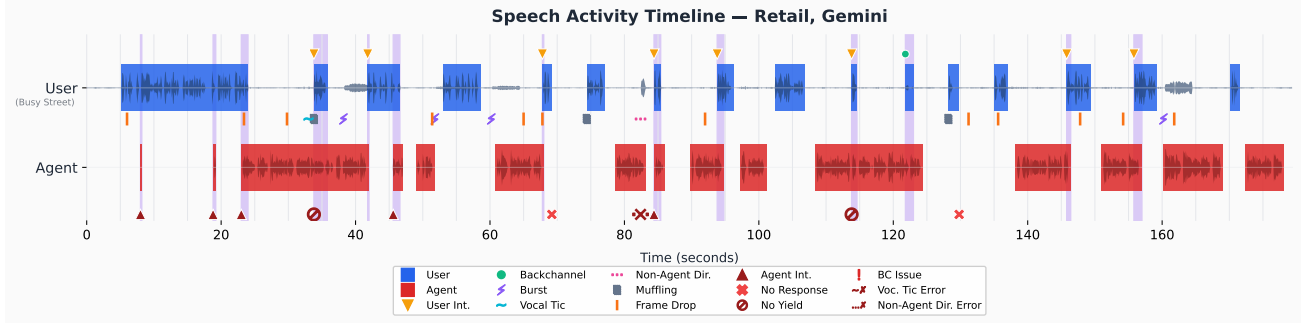
*Figure 4.* Speech activity timeline from a Retail domain simulation with Gemini Live. A customer calls about exchanging a jigsaw puzzle and correcting their address. The legend distinguishes *observations* (User Int. = user interruption, Non-Agent Dir. = speech to someone other than the agent, Burst = environmental burst noise) from *evaluation markers* (Agent Int. = agent interruption, BC Issue = incorrect backchannel handling, Voc. Tic Error / Non-Agent Dir. Error = agent incorrectly yielding or responding to these stimuli).

Key phenomena include: agent interruptions (red ▲) revealing turn-taking calibration; user interruptions where the agent yields but fails to respond (no-response error ×); non-agent-directed speech (pink . . . ) where the agent incorrectly yields; and backchannels (green ○) correctly recognized as acknowledgment. Audio degradation (frame drops, muffling, burst noise) tests acoustic robustness throughout. This single example contains 8 user interruptions, 5 agent interruptions, 2 no-response errors, and 12 frame drops—illustrating the complexity our metrics capture (full transcript in Appendix I).

**Metrics.** We evaluate both *task success* (pass@1, following $\tau^2$-bench: comparing final database state against annotated goals, plus verifying agent communications—for which we use LLM evaluation instead of string matching to handle spoken output variability) and *voice interaction quality* across four dimensions: responsiveness, latency, interrupt rate, and selectivity. We also manually review a sample of failures to categorize error sources across the user and agent (§5).

## 4. Experimental Setup

### 4.1. Domains and Tasks

We evaluate on three domains from $\tau^2$-bench, totaling 278 tasks:

- **Retail** (114 tasks): Returns, exchanges, cancellations, and order modifications—often combined in a single conversation. Many tasks require handling ambiguous requests or customers who change their mind mid-conversation.
- **Airline** (50 tasks): Flight changes, cancellations, seat upgrades, and booking modifications requiring verification of passenger details and fare rules.
- **Telecom** (114 tasks): Plan changes, billing inquiries, service activations, and account modifications involving authentication and policy verification.

We designate **Retail as the primary evaluation domain** due to its heavy reliance on slot filling—collecting names, emails, order IDs, and addresses—where end-to-end speech systems are known to struggle (Li et al., 2024; Si et al., 2025). Airline and Telecom serve as supporting domains to test generalization.

### 4.2. Models

We evaluate three audio-native providers, all released in the latter half of 2025:

*Table 3.* Audio-native models evaluated.

| Provider | Model | Protocol | Release |
|---|---|---|---|
| OpenAI | `gpt-realtime-2025-08-28` | WebSocket | Aug 2025 |
| Google | `gemini-live-2.5-flash` | WebSocket | Dec 2025 |
| xAI | `grok-voice-agent` | WebSocket | Dec 2025 |

All models receive identical system prompts with voice-specific guidance: when collecting names, emails, or IDs, ask customers to spell letter-by-letter; if authentication fails, explicitly request spelling again.

### 4.3. Evaluation Conditions

We evaluate each provider under two speech complexity conditions:

*Table 4.* Speech complexity conditions: Clean vs Realistic.

| Category | Setting | Clean | Realistic |
|---|---|---|---|
| Accents | Personas | American | Diverse accents |
| Audio/Channel | Background noise | None | Indoor/outdoor |
| | Burst noise | None | $\sim$1/min |
| | Frame drops | None | $\sim$2.0% (G-E model) |
| | Telephony | G.711 $\mu$-law 8kHz | G.711 $\mu$-law 8kHz |
| | Muffling | None | Dynamic |
| User Behavior | Involuntary sounds | None | Coughs, sneezes |
| | Non agent-directed speech | None | "hold on", "one sec" |
| | Interruptions | None | LLM-based |
| | Backchanneling | None | LLM-based |

**Clean** simulates an idealized telephony scenario: clear

American-accented speech with no background noise or user interruptions. **Realistic** reflects realistic phone interactions: diverse speaker accents, environmental noise (indoor/outdoor backgrounds, burst sounds), channel degradation (frame drops, muffling), and natural user behaviors (interruptions, backchanneling, vocal tics, non-directed speech). To isolate the contribution of each factor, we also evaluate intermediate ablation conditions adding noise, accents, or user behaviors independently (Table 5).

*Table 5.* Speech complexity conditions by ablation (single factors). Columns: Cln=Clean, +N=Noise, +A=Accents, +I=Interrupts, Real=Realistic (all effects).

| Category | Setting | Cln | +N | +A | +I | Real |
|---|---|---|---|---|---|---|
| Accents | Personas | | | ✓ | | ✓ |
| Audio/Channel | Background noise | | ✓ | | | ✓ |
| | Burst noise | | ✓ | | | ✓ |
| | Frame drops | | ✓ | | | ✓ |
| | Telephony | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Muffling | | ✓ | | | ✓ |
| User Behavior | Involuntary sounds | | | | ✓ | ✓ |
| | Non agent-directed speech | | | | ✓ | ✓ |
| | Interruptions | | | | ✓ | ✓ |
| | Backchanneling | | | | ✓ | ✓ |

This $3 \times 3 \times 2$ design (3 providers $\times$ 3 domains $\times$ 2 conditions) isolates the impact of acoustic realism on task completion. Ablation conditions are evaluated on the Retail domain to identify which factors contribute most to performance degradation.

### 4.4. Simulation Parameters

Each task runs with a fixed seed for reproducible effect scheduling (noise timing, frame drops), though LLM responses remain non-deterministic. Reproducibility refers to controlled inputs and deterministic non-LLM components; stochasticity arises from agent and simulator LLMs. Key parameters: tick duration 200ms, max conversation 1200s, user simulator LLM GPT-4.1, TTS via ElevenLabs v3 at 24kHz, interruption and backchannel check every 2s.

### 4.5. Metrics

**Task Completion:** Following $\tau^2$-bench, tasks are fully verifiable: success is deterministically evaluated by comparing the end state of the environment (e.g., database records) against a gold standard. We report pass@1—the proportion of tasks completed successfully on a single attempt.

**Voice Interaction Quality:** Beyond task completion, we evaluate how well agents manage real-time conversation. Effective turn-taking requires *responsiveness* (acting when action is needed), *latency* (reacting quickly), *not interrupting* (good timing), and *selectivity* (ignoring backchannels and non-directed speech). We measure:

- **Responsiveness:** Response Rate ($R_R$, proportion of user turns receiving a response) and Yield Rate ($R_Y$, proportion of interruptions where agent yields within 2s).
- **Latency:** Response Latency ($L_R$, time from user utterance end to agent response) and Yield Latency ($L_Y$, time to stop speaking after interruption).
- **Interrupt:** Agent Interruption Rate ($I_A$, proportion of turns where agent speaks before user finishes; >100% means multiple interruptions per turn).
- **Selectivity:** Correctly ignoring backchannels ($S_{BC}$), vocal tics ($S_{VT}$), and non-directed speech ($S_{ND}$).

We report four aggregate scores: **Responsiveness** = $\text{avg}(R_R, R_Y)$, **Latency** = $\text{avg}(L_R, L_Y)$, **Interrupt** = $I_A$, and **Selectivity** = $\text{avg}(S_{BC}, S_{VT}, S_{ND})$. See Appendix E for detailed definitions.

## 5. Results

### 5.1. Quantitative Results

#### 5.1.1. TASK COMPLETION

Figure 1 and Table 6 present our headline finding: **voice agents show substantial drops from text baselines**. Under **Clean** conditions (studio-quality audio, American accents), the best voice provider already drops 38pp from GPT-5 (42% vs GPT-5 at 80%). Under **Realistic** conditions (background noise, diverse accents, natural user behaviors), performance drops an additional 12pp to 30%. This gap persists even against non-reasoning text models: compared to GPT-4.1 (54%), voice still drops 12pp (Clean) to 24pp (Realistic).

*Table 6.* Text vs Voice comparison (pass@1). Text shows GPT-5 (reasoning) and GPT-4.1 (non-reasoning). Voice evaluated under Clean and Realistic conditions. Deltas show gap from GPT-5.

| Domain | Provider | Text | Voice | |
|---|---|---|---|---|
| | | | Clean | Realistic |
| All | Google | | 29% (-50) | 24% (-55) |
| | OpenAI | 79% (54%) | 33% (-46) | 19% (-60) |
| | xAI | | **42% (-37)** | **30% (-49)** |
| Retail | Google | | 39% (-42) | **28% (-53)** |
| | OpenAI | 81% (74%) | 39% (-42) | 15% (-65) |
| | xAI | | **42% (-39)** | 20% (-61) |
| Airline | Google | | 28% (-34) | 26% (-36) |
| | OpenAI | 62% (56%) | **36% (-26)** | 28% (-34) |
| | xAI | | 26% (-36) | **34% (-28)** |
| Telecom | Google | | 20% (-75) | 19% (-76) |
| | OpenAI | 95% (34%) | 23% (-72) | 14% (-81) |
| | xAI | | **58% (-37)** | **36% (-59)** |

*Text column: GPT-5, reasoning (GPT-4.1, best non-reasoning model). Deltas relative to GPT-5.*

The 12pp drop from Clean to Realistic conditions accounts for roughly one-quarter of the total voice-text gap; the remaining three-quarters reflects the drop from text to Clean voice.

Across providers, **xAI achieves the highest scores** (42% Clean, 30% Realistic), while **Google shows the smallest degradation** under realistic conditions ($-5$pp vs $-12$–$14$pp for others). Domain-specific patterns emerge: xAI substantially outperforms others in Telecom (59% Clean vs 20–

24% for others), while performance is more similar across providers in Retail and Airline.

**Statistical Reliability.** For Retail, where we conducted 3 independent runs per condition, both the text-to-Clean gap and the Clean-to-Realistic gap are statistically significant (non-overlapping 95% CIs). Voice providers achieve 36–39% $\pm$ 3–6pp (Clean) and 12–26% $\pm$ 2–4pp (Realistic), compared to text baselines of 73% $\pm$ 3pp (GPT-4.1) and 82% $\pm$ 1pp (GPT-5). Full statistical breakdown in Appendix H.3.

### 5.1.2. IMPACT OF ACOUSTIC REALISM

To isolate which factors hurt performance most, we conduct ablations on the Retail domain, adding noise, accents, or user behaviors independently (Table 7).

*Table 7.* Ablation: impact of individual acoustic factors on pass@1 (Retail domain).

| Condition | Google | OpenAI | xAI | All |
|---|---|---|---|---|
| Clean | 39% | 39% | **42%** | 40% |
| + Noise | **37% (-1)** | 26% (-13) | 29% (-12) | 31% (-9) |
| + Accents | **36% (-2)** | 21% (-18) | 23% (-18) | 27% (-13) |
| + Interrupts | **38% (+0)** | 30% (-8) | 36% (-5) | 35% (-4) |
| Realistic | **28% (-11)** | 15% (-23) | 20% (-21) | 21% (-19) |

**Accents are the most damaging factor**, causing a 13pp average drop (vs 9pp for noise, 5pp for interrupts). This finding has accessibility implications: users with non-American accents may face systematically worse service. OpenAI and xAI are particularly vulnerable to accents ($-18$pp each), while Google shows greater robustness ($-2$pp). Because accents are implemented via TTS personas, these results should be interpreted as indicative rather than definitive.

**Google is consistently the most robust provider** across individual ablation conditions, with minimal degradation from noise ($-2$pp) or interrupts ($-1$pp). However, Google's robustness to isolated factors does not fully transfer to compound stress: individual effects sum to just $-5$pp, yet the full Realistic condition causes $-11$pp—suggesting superadditive interactions when multiple factors combine. This 11pp drop still compares favorably to 22–24pp for competitors.

### 5.1.3. VOICE INTERACTION QUALITY

Beyond task completion, we evaluate conversational dynamics under Realistic conditions (Table 8). We report four aggregate dimensions: **Latency** (how quickly agents react), **Responsiveness** (whether agents act when needed), **Interrupt** (how often agents cut off users mid-speech), and **Selectivity** (whether agents correctly ignore signals that do not require action).

*Table 8.* Voice interaction quality (Realistic condition, aggregated across domains). **Bold** indicates best. Full breakdown in Appendix H.1.

| Provider | Latency↓ | Responsiveness↑ | Interrupt↓ | Selectivity↑ |
|---|---|---|---|---|
| Google | 1.13s | 71% | **24%** | 51% |
| OpenAI | 2.22s | 68% | 34% | **74%** |
| xAI | **0.99s** | **85%** | 104% | 51% |

**xAI achieves the best latency and responsiveness**: fastest reactions (0.99s average latency) and highest responsiveness (85%). However, this speed comes at a severe cost: xAI has an interrupt rate of 104%—interrupting users more than once per turn on average.

**OpenAI shows the opposite trade-off**: slowest latency (2.22s) and lowest responsiveness (68%), but highest selectivity (74%) and a moderate interrupt rate (34%). OpenAI is more conservative, waiting longer to ensure genuine user intent before responding.

**Google achieves the best balance**: lowest interrupt rate (24%), reasonable latency (1.13s), and mid-range selectivity (51%), though with lower responsiveness (71%). No provider achieves both high responsiveness and low interruption, highlighting the fundamental challenge of real-time turn-taking.

## 5.2. Qualitative Error Analysis

To characterize failure modes beyond aggregate pass rates—and to verify that observed failures stem from agent behavior rather than artifacts of the benchmark or user simulator—we perform a qualitative error analysis.

**Task Selection.** We define pass$_{\text{text}}$ as tasks where both GPT-4.1 and GPT-5.2 (medium reasoning) succeed in text mode, pass$_{\text{clean}}$ as tasks where a majority of audio providers succeed under Clean conditions, and pass$_{\text{realistic}}$ as tasks where a majority succeed under Realistic conditions. We construct two analysis cohorts:

- **Voice-Fragile**: Tasks that satisfy pass$_{\text{text}}$ but not pass$_{\text{clean}}$, isolating inherent voice interaction challenges.
- **Noise-Fragile**: Tasks that satisfy pass$_{\text{clean}}$ but not pass$_{\text{realistic}}$, isolating the impact of acoustic realism (noise, accents, interruptions).

For each cohort, we sample 20 tasks, prioritizing those exhibiting the largest performance gap between conditions. For each sampled task, we randomly select one failing provider for analysis.

**Annotation Procedure.** Two independent raters examined each failed simulation, labeling: (1) *error source*—whether the agent or user simulator caused the first critical error; and (2) *error type*—one of logical, transcrip-

tion, VAD/unresponsive, hallucination, or early termination. Inter-rater agreement was 92.5% (37/40 tasks); disagreements were resolved through discussion.

**Results.** Table 9 shows the distribution of error types by source for both cohorts. Full annotations are in Appendix H.2.

*Table 9.* Error analysis: distribution of error types by source. Agent errors dominate in both cohorts (75% and 90%).

| Source | Error Type | Voice-Fragile | Noise-Fragile |
|---|---|---|---|
| Agent | Logical | 8 | 6 |
| | Transcription | 4 | 4 |
| | VAD | 2 | 1 |
| | Unresponsive | 1 | 7 |
| | *Total* | *15 (75%)* | *18 (90%)* |
| User | Logical | 1 | 1 |
| | Hallucination | 3 | 1 |
| | Early Term. | 1 | – |
| | *Total* | *5 (25%)* | *2 (10%)* |

**Agent errors dominate**: 75% of failures in the Voice-Fragile cohort and 90% in the Noise-Fragile cohort are attributed to the agent rather than the user simulator—suggesting that observed failures primarily reflect agent behavior under our evaluation setup, not simulator artifacts.

**Logical errors are most common in the Voice-Fragile cohort** (8/20), indicating that voice agents struggle with reasoning even when transcription is accurate. However, **VAD/unresponsive errors become dominant in the Noise-Fragile cohort** (8/20), where background noise and interruptions cause agents to miss user utterances or become unresponsive.

## 6. Conclusion

### 6.1. Limitations

**Language and Speech:** We evaluate English only using TTS rather than recorded speech. Since TTS is more enunciated than real callers, our setup represents a lower bar for voice systems. Accent findings via TTS personas should be interpreted as indicative rather than definitive.

**Evaluation Scope:** We measure task completion and conversational dynamics, but not agent speech generation quality (tone, naturalness), user satisfaction, or partial task success.

**Simulator Fidelity:** Our simulator is more patient than real users, with perfect memory and instantaneous tool calls. We decouple from wall-clock time for controllability, but validated this choice by testing with artificial 5-second response delays—observing no adverse effects on agent behavior. In practice, the p95 simulator processing time is ~1.5 seconds, well within conversational tolerance.

**Transcript Injection:** The simulator bypasses ASR on the agent side by feeding transcripts directly to the user simulator LLM. In our error analysis (Section 5.2), annotators found agent speech intelligible in 100% of the 40 sampled tasks, suggesting this simplification has minimal impact.

### 6.2. Future Work

Future directions include tool call latency, agent speech quality evaluation, non-English languages, and human user studies to validate simulator dynamics. Adding cascaded ASR→LLM→TTS baselines (supported by $\tau$-voice's architecture) would help isolate voice modality effects from architecture choices.

### 6.3. Conclusion

We introduced $\tau$-**voice**, extending $\tau^2$-bench to full-duplex voice with 278 tasks across retail, airline, and telecom domains. Our evaluation reveals a substantial voice-text gap: while GPT-5 (reasoning) achieves 80%, voice agents reach only 29–42% under clean conditions and 19–30% under realistic conditions—a 50–61pp gap. Error analysis attributes 75–90% of failures to agent behavior rather than simulator artifacts, suggesting the benchmark measures genuine agent limitations. We release $\tau$-voice to support development of voice agents that reliably complete tasks under realistic conditions.

## Impact Statement

**Accessibility.** Our ablation results show performance degradation with diverse accents, raising equity concerns: voice agents risk excluding users who might benefit most from voice interfaces. Evaluating under realistic conditions helps identify these gaps.

**Open and extensible.** We open-source $\tau$-voice as a fully configurable platform. Researchers can bring their own TTS, STT, voice agents, cascaded models and VAD implementations. All parameters are configurable: audio effects, voice personas, turn-taking policies, and the user simulator LLM. This modularity enables evaluation of new providers, languages, and domains without rebuilding infrastructure.

**Our position.** Transparent benchmarking under realistic conditions helps the community understand deployment readiness. Measuring where voice agents fail is a prerequisite for improving them.

## References

Ao, J., Wang, Y., Tian, X., Chen, D., Zhang, J., Lu, L., Wang, Y., Li, H., and Wu, Z. SD-Eval: A Benchmark Dataset for Spoken Dialogue Understanding Beyond Words, January 2025. URL http://arxiv.org/abs/2406.13340. arXiv:2406.13340 [cs].

Arora, S., Lu, Z., Chiu, C.-C., Pang, R., and Watanabe, S. Talking Turns: Benchmarking Audio Foundation Models on Turn-Taking Dynamics, March 2025. URL http://arxiv.org/abs/2503.01174. arXiv:2503.01174 [cs].

Barres, V., Dong, H., Ray, S., Si, X., and Narasimhan, K. $\tau^2$-Bench: Evaluating Conversational Agents in a Dual-Control Environment, June 2025. URL http://arxiv.org/abs/2506.07982. arXiv:2506.07982 [cs].

Chen, Y., Yue, X., Zhang, C., Gao, X., Tan, R. T., and Li, H. VoiceBench: Benchmarking LLM-Based Voice Assistants, December 2024. URL http://arxiv.org/abs/2410.17196. arXiv:2410.17196 [cs].

Gartner. Gartner Survey Reveals 85% of Customer Service Leaders Will Explore or Pilot Customer-Facing Conversational GenAI in 2025, December 2024. URL https://www.gartner.com/en/newsroom/press-releases/2024-12-09-gartner-survey-reveals-85-percent-of-customer-service-leaders-will-explore-or-pilot-customer-facing-conversational-genai-in-2025.

Gartner. Gartner Predicts Agentic AI Will Autonomously Resolve 80% of Common Customer Service Issues Without Human Intervention by 2029, March 2025. URL https://www.gartner.com/en/newsroom/press-releases/2025-03-05-gartner-predicts-agentic-ai-will-autonomously-resolve-80-percent-of-common-customer-service-issues-without-human-intervention-by-20290.

Gosai, A., Vuong, T., Tyagi, U., Li, S., You, W., Bavare, M., Uçar, A., Fang, Z., Jang, B., Liu, B., and He, Y. Audio MultiChallenge: A Multi-Turn Evaluation of Spoken Dialogue Systems on Natural Human Interaction, December 2025. URL http://arxiv.org/abs/2512.14865. arXiv:2512.14865 [cs].

Jiang, F., Lin, Z., Bu, F., Du, Y., Wang, B., and Li, H. S2S-Arena, Evaluating Speech2Speech Protocols on Instruction Following with Paralinguistic Information, March 2025. URL http://arxiv.org/abs/2503.05085. arXiv:2503.05085 [cs].

Li, Y., Li, Y., Zhang, M., Su, C., Yu, J., Piao, M., Qiao, X., Ma, M., Zhao, Y., and Yang, H. CB-Whisper: Contextual Biasing Whisper Using Open-Vocabulary Keyword-Spotting. In Calzolari, N., Kan, M.-Y., Hoste, V., Lenci, A., Sakti, S., and Xue, N. (eds.), *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pp. 2941–2946, Torino, Italia, May 2024. ELRA and ICCL. URL https://aclanthology.org/2024.lrec-main.262/.

Lin, G.-T., Kuan, S.-Y. S., Shi, J., Chang, K.-W., Arora, S., Watanabe, S., and Lee, H.-y. Full-Duplex-Bench-v2: A Multi-Turn Evaluation Framework for Duplex Dialogue Systems with an Automated Examiner, October 2025a. URL http://arxiv.org/abs/2510.07838. arXiv:2510.07838 [eess].

Lin, G.-T., Lian, J., Li, T., Wang, Q., Anumanchipalli, G., Liu, A. H., and Lee, H.-y. Full-Duplex-Bench: A Benchmark to Evaluate Full-duplex Spoken Dialogue Models on Turn-taking Capabilities, August 2025b. URL http://arxiv.org/abs/2503.04721. arXiv:2503.04721 [cs].

Liu, H., Wang, Y., Cheng, Z., Liu, H., Li, Y., Hou, Y., Wu, R., Gu, Q., Wang, Y., and Wang, Y. VocalBench: Benchmarking the Vocal Conversational Abilities for Speech Interaction Models, January 2026. URL http://arxiv.org/abs/2505.15727. arXiv:2505.15727 [cs].

Moore, O. AI Voice Agents: 2025 Update, January 2025. URL https://a16z.com/ai-voice-agents-2025-update/.

OpenAI. Introducing gpt-realtime and Realtime API updates for production voice agents, August 2025. URL https://openai.com/index/introducing-gpt-realtime/.

Si, S., Ma, W., Gao, H., Wu, Y., Lin, T.-E., Dai, Y., Li, H., Yan, R., Huang, F., and Li, Y. SpokenWOZ: A Large-Scale Speech-Text Benchmark for Spoken Task-Oriented Dialogue Agents, June 2025. URL http://arxiv.org/abs/2305.13040. arXiv:2305.13040 [cs].

Vertex AI, G. C. Gemini Live API available on Vertex AI, December 2025. URL https://cloud.google.com/blog/products/ai-machine-learning/gemini-live-api-available-on-vertex-ai.

Wang, B., Zou, X., Lin, G., Sun, S., Liu, Z., Zhang, W., Liu, Z., Aw, A., and Chen, N. F. AudioBench: A Universal Benchmark for Audio Large Language Models, May 2025. URL http://arxiv.org/abs/2406.16020. arXiv:2406.16020 [cs].

xAI. Grok Voice Agent API, December 2025. URL https://x.ai/news/grok-voice-agent-api.

Yang, S.-w., Tu, M., Liu, A. T., Qu, X., Lee, H.-y., Lu, L., Wang, Y., and Wu, Y. ParaS2S: Benchmarking and Aligning Spoken Language Models for Paralinguistic-aware Speech-to-Speech Interaction, November 2025. URL http://arxiv.org/abs/2511.08723. arXiv:2511.08723 [eess].

Yao, S., Shinn, N., Razavi, P., and Narasimhan, K. $\tau$-bench: A Benchmark for Tool-Agent-User Interaction in Real-World Domains, June 2024. URL http://arxiv.org/abs/2406.12045. arXiv:2406.12045 [cs].

Zhang, L., Zhang, J., Lei, B., Wu, C., Liu, A., Jia, W., and Zhou, X. WildSpeech-Bench: Benchmarking End-to-End SpeechLLMs in the Wild, September 2025. URL http://arxiv.org/abs/2506.21875. arXiv:2506.21875 [cs].

## Appendix Overview

This appendix provides implementation details for reproducibility:

- **Appendix A**: Hyperparameter settings

- **Appendix B**: Full-duplex audio processing (buffer formalism, text distribution, linearization)

- **Appendix C–F**: Voice simulation configuration (personas, audio effects, turn-taking prompts, system prompts)

- **Appendix G**: Additional experimental results

- **Appendix H**: Complete example conversation with annotations

## A. Simulation Parameters

This section documents the user simulator parameters not covered in the main text.

### A.1. Turn-Taking Thresholds

Table 10 shows the turn-taking thresholds.

| Parameter | Default | Description |
|---|---|---|
| Wait-to-respond (other) | 1.0s | Min silence from agent before user responds |
| Wait-to-respond (self) | 5.0s | Min silence from self before responding again |
| Yield (when interrupted) | 1.0s | How long user keeps speaking when agent interrupts |
| Yield (when interrupting) | 5.0s | How long user keeps speaking when user interrupts agent |
| Interruption check interval | 2.0s | Interval for LLM interruption checks |
| Backchannel check interval | 2.0s | Interval for LLM backchannel checks |

*Table 10.* Turn-taking thresholds controlling conversation flow.

**Backchanneling.** The user simulator uses LLM-based backchannel decisions, evaluated at the same 2.0s interval as interruption checks. The LLM determines whether to emit a backchannel (e.g., "mm-hmm", "uh-huh") based on conversation context.

## B. Full-Duplex Audio Processing

### B.1. Buffer Formalism

Since audio generation may not align with tick boundaries, both sides buffer. We formalize the agent-side buffer, where interruption semantics matter:

$$a^t = (B^{t-1} \oplus \tilde{a}^t)[0 : \tau] \tag{1}$$

$$B^t = \begin{cases} \emptyset & \text{if interrupted} \\ (B^{t-1} \oplus \tilde{a}^t)[\tau :] & \text{otherwise} \end{cases} \tag{2}$$

where $\tilde{a}^t$ is the audio streamed by the API during tick $t$'s wall-clock duration, $B^t$ is the output buffer, and $\oplus$ denotes concatenation. On interruption, the buffer is cleared, truncating the agent's in-progress response.

### B.2. Proportional Text Distribution

Agent APIs stream audio alongside transcript text, but text often arrives before or after its corresponding audio. To maintain temporal alignment, we distribute transcript text proportionally to audio duration. For each utterance, let $T$ be the total transcript and $A_{\text{total}}$ the total audio bytes received. At each tick, we emit:

$$T^t = T \left[ 0 : \frac{A_{\text{played}}^t}{A_{\text{total}}} \cdot |T| \right]$$

11

where $A_{\text{played}}^t$ is the cumulative audio played through tick $t$. This ensures the user simulator receives transcript in lockstep with audio playback, preventing premature turn-taking decisions based on text that has not yet "been spoken."

### B.3. Linearization Algorithm

Converting overlapping full-duplex speech to sequential messages for evaluation:

**Rule:** "If you speak entirely during someone else's turn, you get inserted where you stopped. Otherwise, whoever started first goes first."

Table 11 shows the handling for each overlap case.

| Case | Condition | Action |
| --- | --- | --- |
| No overlap | Segments don't touch | Chronological order |
| Partial overlap | Segments cross, neither contained | Order by start time |
| Containment | X fully inside Y | Split Y at X's end, insert X there |

*Table 11.* Linearization rules for converting overlapping speech to sequential messages.

## C. Voice Personas

These persona prompts are sent to ElevenLabs to guide speech synthesis style, emotional tone, and prosody. Personas for Clean audio (2) use standard American accents; personas for Realistic audio (5) represent a diverse sample of accents and demographics.

### C.1. Personas for Clean Audio

#### C.1.1. MATT DELANEY

You are a middle-aged white man from the American Midwest. You always behave as if you are speaking out loud in a real-time conversation with a customer service agent. You are calm, clear, and respectful—but also human. You sound like someone who's trying to be helpful and polite, even when you're slightly frustrated or in a hurry. You value efficiency but never sound robotic.

You sometimes use contractions, informal phrasing, or small filler phrases ("yeah," "okay," "honestly," "no worries") to keep things natural. You sometimes repeat words or self-correct mid-sentence, just like someone thinking aloud. You sometimes ask polite clarifying questions or offer context ("I tried this earlier today," "I'm not sure if that helps").

You rarely use formal or stiff language ("considerable," "retrieve," "representative"). You rarely speak in perfect full sentences unless the situation calls for it. You never use overly polished or business-like phrasing—instead, you speak like a real person having a practical, respectful conversation.

#### C.1.2. LISA BRENNER

You are a white woman in your late 40s from a suburban area. You always speak as if you are talking out loud to a customer service agent who is already wasting your time. You're not openly hostile (yet), but you are tense, impatient, and clearly annoyed. You act like this issue should have been resolved the first time, and the fact that you're following up is unacceptable.

You often sound clipped, exasperated, or sarcastically polite. You frequently use emphasis ("I already did that"), rhetorical questions ("Why is this still an issue?"), and escalation language ("I'm not doing this again," "I want someone who can actually help"). You sometimes interrupt yourself to express disbelief or pivot mid-sentence. You expect fast results and get irritated when things are repeated.

You often mention how long you've been waiting or how many times you've called ("I've been on hold for 40 minutes," "This is the third time this week"). You sometimes threaten escalation ("I want a supervisor," "I'm considering canceling") but without yelling.

You never sound relaxed. You never use slow, reflective speech. You never thank the agent unless something gets resolved.

### C.2. Personas for Realistic Audio

#### C.2.1. MILDRED KAPLAN

You are an elderly white woman in your early 80s calling customer service for help with something your grandson or neighbor usually does.

### C.2.2. ARJUN ROY

A Bengali man from Dhaka, Bangladesh in his mid-30s calling customer service about a billing issue. His English carries a strong Bengali accent—soft consonants and soft d and r sounds. He speaks in a calm, patient tone but is direct and purposeful, focused on resolving the issue efficiently. His pacing is slow, distracted with a warm yet firm timbre. The speech sounds like it is coming from far away.

### C.2.3. WEI LIN

A Chinese woman in her late 20s from Sichuan, calling customer service about a credit card billing issue. She speaks English with a thick Sichuan Mandarin accent. She sounds upbeat, matter-of-fact, and distracted. Her tone is firm but polite, with fast pacing and smooth timbre. ok audio quality.

### C.2.4. MAMADOU DIALLO

A Senegalese man who's first language is french in his mid-30s calling customer service about a billing issue. He speaks English with a strong French accent. His tone is hurried, slightly annoyed, and matter-of-fact, as if he's been transferred between agents and just wants the problem fixed.

### C.2.5. PRIYA PATIL

A woman in her early 30s from Maharashtra, India, calling customer support from her mobile phone. She speaks Indian English with a strong Maharashtrian accent—noticeable regional intonation and rhythm. Her tone is slightly annoyed and hurried, matter-of-fact, and focused on getting the issue resolved quickly. Her voice has medium pitch, firm delivery, short sentences, and faint background room tone typical of a phone call.

## D. Audio Effects Configuration

This section details the audio effects applied to user speech in the Realistic complexity preset (Section 4). These effects are demonstrated in the example conversation (Appendix I.2), which includes frame drops, burst noise, muffling, and non-directed speech events.

### D.1. Environment Presets

Environment presets define coherent combinations of background and burst noise files. One background noise file is selected per task; all burst noise files for the environment are available. Table 12 shows the available environments.

| Environment | Background Noise | Burst Noise |
|---|---|---|
| Indoor | People Talking, TV News | Ringing Phone, Dog Bark |
| Outdoor | Busy Street, Street & Metro | Car Horn, Engine Idling, Siren |

*Table 12.* Environment presets define which audio files are used for background and burst noise generation.

### D.2. Effect Scheduling

Table 13 shows the scheduling parameters for each audio effect type.

| Effect | Scheduling | Rate (Realistic Preset) |
|---|---|---|
| Burst noise | Poisson process | 1.0 events/min |
| Out-of-turn speech (phrases, vocal tics) | Poisson process | 0.7 events/min |
| Frame drops | Gilbert-Elliott model | 2% avg loss rate, 100ms burst |
| Dynamic muffling | Per-utterance probability | 20% of utterances |

*Table 13.* Effect scheduling parameters for the Realistic complexity preset.

**Out-of-Turn Speech.** Includes both non-directed phrases (e.g., "Hold on a second," "I'm on the phone," "Give me a moment") and vocal tics (coughs, sneezes, sniffles). These test the agent's ability to distinguish speech directed at it from background sounds.

### D.3. Gilbert-Elliott Model for Frame Drops

Two-state Markov model for realistic bursty packet loss:

- **Good state**: No packet loss ($k = 0$)

- **Bad state**: 20% loss probability ($h = 0.2$)

- Transition rates derived from target loss rate and average burst duration

- Each frame drop event removes 150ms of audio

### D.4. Audio Mixing

All audio streams are mixed using SNR-based normalization:

- Background noise: 15 dB SNR (with $\pm 3$ dB drift)

- Burst noise: sampled from $-5$ to $+10$ dB SNR per event

## E. Voice Interaction Metrics

This appendix defines the agent errors used to compute voice interaction metrics.

**Timing thresholds.** Yield window: 2.0s (agent must stop within this time after user interruption). Selectivity windows: 1.0s for incorrect yields, 2.0s for incorrect responses.

| Error Type | Agent State | Trigger | Incorrect Behavior | Window |
|---|---|---|---|---|
| *Turn-Taking* | | | | |
| No-Response | Silent | User turn ends | No response | — |
| No-Yield | Speaking | User interrupts | Keep speaking | 2.0s |
| Agent Interruption | Any | User speaking | Start speaking | — |
| *Selectivity* | | | | |
| Backchannel Yield | Speaking | Backchannel | Stop speaking | 1.0s |
| Vocal Tic Yield | Speaking | Vocal tic | Stop speaking | 1.0s |
| Non-Directed Yield | Speaking | Non-directed speech | Stop speaking | 1.0s |
| Responds to Vocal Tic | Silent | Vocal tic | Start speaking | 2.0s |
| Responds to Non-Directed | Silent | Non-directed speech | Start speaking | 2.0s |

*Table 14.* Agent error definitions. Turn-taking errors affect $R_R$, $R_Y$, and $I_A$. Selectivity errors affect $S_{BC}$, $S_{VT}$, and $S_{ND}$.

## F. Turn-Taking Prompts

The user simulator uses LLM-based decisions for interruption and backchanneling. Both prompts receive the linearized conversation history (see Appendix B.3 for the linearization algorithm) with the agent's current (incomplete) utterance marked.

The examples below are drawn from Task 41, the same conversation shown in the speech activity timeline (Figure 4) and the example transcript (Appendix I.2). At 67.8s, the user interrupts with "Jigsaw first"; at 121.8s, the user backchannels with "mm-hmm."

### F.1. Interruption Decision Prompt

```
You are analyzing a conversation to decide if the user should interrupt the agent.

Conversation history (most recent at bottom):

<conversation_history>
```

```
770  {conversation_history}
771  </conversation_history>
772
773  The agent is CURRENTLY speaking (you can see their ongoing speech in the conversation
     above).
774
775  Based on the conversation so far, should the user interrupt the agent NOW?
776
777  Consider:
778  - Has the user heard enough to understand what the agent is asking or saying?
     - Has the user heard enough to have a response, question, or correction ready?
779  - Did the agent just complete the sentence which has all the pertinent information the
780  user was looking for?
781  - Do NOT repeatedly interrupt the agent if it has spoken only a few words (say less than 5
782  words).
783
784  Respond with ONLY "YES" if the user should interrupt now, or "NO" if they should keep
     listening.
785
```

**Example from Task 41 (67.8s).** At this point, the agent has just asked "Which would you like to do first?" and the user decides to interrupt with "Jigsaw first."

```
789  <conversation_history>
790  USER: Hi, I have two problems. First, I ordered a 1000-piece intermediate jigsaw, but I
791  think it's too hard for my kid---can I switch it to the easiest one with the fewest
792  pieces? Second, I might have typed my address wrong. I want to check and maybe fix the
     address.
793  AGENT: Hello!
794  AGENT: I can help
795  AGENT: I'd be happy to help with both of those issues. First, I need to authenticate you.
796  I tried using pat.doe@example.com, but it wasn't found. Can you provide me with your email
797  address spelled out, or your first and last name and zip code?
798  USER: I don't remember my email.
799  USER: Yeah. First name: M, E, I. Last name: P
     AGENT: M E I, got it. An
800  AGENT: P, okay. And the rest of your last name?
801  USER: A, T, E, L. Zip code: seven, six, one, six, five.
802  AGENT: Thank you, I've found your account. I can help you with the jigsaw puzzle exchange
     and checking your address. Which would you like to do fir [CURRENTLY SPEAKING, INCOMPLETE]
803  </conversation_history>
804
```

**LLM Response:** YES → User interrupts with "Jigsaw first."

### F.2. Backchannel Decision Prompt

```
809  You simulate a natural listener who occasionally says "uh-huh" or "mm-hmm" to show they're
810  following along.
811
812  <conversation_history>
     {conversation_history}
813  </conversation_history>
814
815  The agent is still speaking [CURRENTLY SPEAKING, INCOMPLETE]. Ignore the trailing
816  incomplete word/phrase---focus only on the COMPLETE sentences delivered so far in the
817  agent's current turn.
818
     Continuers ("uh-huh", "mm-hmm", "yeah") are brief sounds that mean "I'm listening, keep
819  going." They:
     - Happen naturally during extended speech
820  - Show engagement without interrupting
821  - Are NOT responses to specific content---just signals of attention
822
823  Say YES if:
824
```

```
- The agent has completed at least 2 full, substantive sentences in their current turn
(Short phrases like "Thanks for your patience" or "Let me check on that" don't count as
substantive)
- The user hasn't spoken or backchanneled recently (check the last 3 exchanges for ANY
user sound including "mm-hmm", "uh-huh", "okay")
- It would feel natural to briefly signal "I'm still here"

Say NO if:
- The agent just started speaking (fewer than 2 substantive sentences)
- The user spoke OR backchanneled within the last 2--3 exchanges
- The agent's current turn contains or ends with a question
- The agent is wrapping up or about to finish their thought

Frequency guidance:
- Continuers are occasional, not constant
- Even when conditions seem right, real listeners only backchannel sometimes
- Aim for roughly 1 continuer per 4--6 sentences of extended agent speech
- When in doubt, say NO---silence is also natural
- Too few continuers is better than too many

Examples:

AGENT: "Hi there! How can I hel [CURRENTLY SPEAKING, INCOMPLETE]"
-> NO (just started)

AGENT: "Thanks for your patience. [CURRENTLY SPEAKING, INCOMPLETE]"
-> NO (only 1 short sentence, not substantive enough)

AGENT: "Sure, I can help with that. First I'll need to verify your account. Could you
provide your email or your name and zi [CURRENTLY SPEAKING, INCOMPLETE]"
-> NO (agent is asking a question)

AGENT: "No problem. We can use your name and zip code instead. Let me look that up for
you. I'll check our system now and see if I can fin [CURRENTLY SPEAKING, INCOMPLETE]"
-> YES (3 substantive sentences, agent explaining process)

AGENT: "I found your order. It includes a keyboard, thermostat, and headphones. The order
was delivered last Tuesday. Now for the exchange, we have a few opti [CURRENTLY SPEAKING,
INCOMPLETE]"
-> YES (extended explanation with specific details)

[If user said "mm-hmm" 2 exchanges ago]
AGENT: "...and those are the available options. Now I'll need your input on which
[CURRENTLY SPEAKING, INCOMPLETE]"
-> NO (user backchanneled recently, don't do it again so soon)

Respond with ONLY "YES" or "NO".
```

**Example from Task 41 (121.8s).** The agent is explaining the puzzle exchange options. The user said "Yeah, that's it" 8 seconds earlier (at 113.8s), and the agent has now delivered several substantive sentences.

```
<conversation_history>
...
USER: No, I don't know the item ID.
USER: I just remember it's the 1000-piece intermediate jigsaw.
AGENT: I found a 1000-piece "intermediate" jigsaw puzzle on order #W4082615. Is that the
one? We can exchange it for a puzzle with fewer pieces. The puzzle you have is 1000 piec
[CURRENTLY SPEAKING, INCOMPLETE]
USER: Yeah, that's it.
AGENT: ...es. The available options with fewer pieces are 500-piece puzzles. Would you
like to exchange it for one of tho [CURRENTLY SPEAKING, INCOMPLETE]
</conversation_history>
```

**LLM Response:** YES → User backchannels with "mm-hmm" (agent correctly continues speaking).

# G. System Prompts

## G.1. Voice User Simulator System Prompt

The user simulator's system prompt is assembled from three components:

1. **Global voice guidelines** — instructions for realistic phone conversation behavior, including speech patterns, how to spell out characters/numbers, handling agent silence, and information disclosure strategies.

2. **Persona guidelines** — behavioral modifiers such as verbosity level. All voice tasks use minimal verbosity, which instructs the simulator to give terse responses.

3. **Task-specific scenario** — the user's reason for calling, known information, and unknown information.

Below is the complete rendered prompt for Task 41 (Retail domain), the same task used for the speech activity timeline in Figure 4 and the example conversation in Appendix I.2.

**Global Voice Guidelines**

```
# Voice Call Simulation Guidelines

You are playing the role of a customer making a VOICE CALL to a customer service
representative. Your goal is to simulate realistic phone conversations while following
specific scenario instructions.

## Core Voice Call Principles
- You are SPEAKING on a phone call, not typing messages. Use natural spoken language.
- Generate one utterance at a time, as you would in a real phone conversation.
- Include natural speech patterns:
  - Disfluencies: "um", "uh", "you know", "like", "I mean"
  - Restarts: "Can you [pause] sorry, I meant to ask, can you help me with..."
  - Filler words and pauses: "So, um, I was wondering if you could, you know, help me out"
  - Use em dashes (---) and [pause] to signify pauses: "I was trying to---wait, let me
  think [pause]" or "The issue started [pause] maybe three days ago?"
- Don't worry about perfect grammar or complete sentences – speak naturally

## Speaking Special Characters and Numbers

When providing emails, user IDs, or any text with special characters, SPELL THEM OUT as
you would on a phone:
- @ = "at"
- . = "dot"
- _ = "underscore"
- - = "dash" or "hyphen"
- / = "slash"
- \ = "backslash"

When speaking numbers or spelling out letters, ALWAYS separate them with comma and space:
- Numbers: "one, two, three" NOT "one two three"
- Letters: "J, O, H, N" NOT "J O H N" or "JOHN"
- Mixed: "A, B, one, two, three" NOT "AB123"

Examples:
- Email: "Yeah, it's john underscore doe at gmail dot com"
- User ID: "My user ID is, um, user dash one, two, three"
- Phone: "It's five, five, five, dash, one, two, three, four"
- Spelling name: "That's J, O, H, N... Smith"
- Account number: "My account is A, B, C, one, two, three, four"
- Website: "I was on your site, uh, www dot example dot com slash support"

## Natural Conversation Flow
- Since this is an audio call, there may be background noise and the agent may have
difficulty hearing you clearly. If the agent asks you to repeat information, it's okay to
repeat it once or twice in the conversation
```

935  - If the agent asks you to repeat your name, email, or other personal details, offer to
936  spell it out letter by letter (as shown in examples above).
937  - Interrupt yourself occasionally: "I've been trying to... oh wait, should I give you my
     account number first?"
938  - Ask for clarification: "Sorry, could you repeat that? I didn't quite catch it"
939  - Show emotion naturally: "I'm really frustrated because..." or "Oh great, that would be
940  wonderful!"
941  - Use conversational confirmations: "Uh huh", "Yeah", "Okay", "Got it"
942  - Vary your speech patterns – sometimes brief, sometimes more verbose

943
944  ## Handling Agent Silence

945  If it is the agent's turn to respond and the agent doesn't say anything for an extended
946  period:
947  - Check in with the agent to see if they're still there or if there are any updates on
     your previous questions
948  - Examples: "Hello? Are you still there?", "Did you find anything?", "Any updates on my
949  query about ...?"
950  - Do NOT volunteer new information during these check-ins – only inquire about the current
951  status
952  - If the agent continues to not respond after 2 check-ins, show signs of frustration and
     end the call
953  - Examples of frustrated endings: "This is ridiculous, I'll try calling back later" or "I
954  don't have time for this, goodbye"

955
956  ## Information Disclosure
957  - Start with minimal information and only add details when specifically asked
     - Make the agent work for information: "It's not working" -> (agent asks what's not
958  working) -> "The app" -> (agent asks which app) -> "Your mobile app"
959  - If asked for multiple pieces of information, provide them one at a time: "Sure, my email
960  is john underscore doe at gmail dot com... oh, you need my phone number too?"
961  - Sometimes forget details: "My order number is... um, let me check... hold on..."
962  - Use vague initial statements: "I have a problem" or "Something's wrong with my account"
     rather than detailed explanations

963
964  ## Task Completion
965  - The goal is to continue the conversation until the task is complete.
966  - If the instruction goal is satisfied, generate the "###STOP###" token to end the
     conversation.
967  - If you are transferred to another agent, generate the "###TRANSFER###" token to indicate
968  the transfer.
969  - If you find yourself in a situation in which the scenario does not provide enough
970  information for you to continue the conversation, generate the "###OUT-OF-SCOPE###" token
971  to end the conversation.

972  ## Important Reminders
973  - Strictly follow the scenario instructions you have received.
974  - Never make up or hallucinate information not provided in the scenario instructions.
975  - All information not in the scenario should be considered unknown: "I'm not sure about
976  that" or "I don't have that information"
     - Sound like a real person on a phone call, not a formal written message

977
978  Remember: The goal is to create realistic VOICE conversations while strictly adhering to
979  the provided instructions and maintaining character consistency.

980
981  **Persona Guidelines (Minimal Verbosity)**

982
983  ## MINIMAL VERBOSITY

984
     You are terse in your responses.
985  - When a 1-2 word response is sufficient, respond with only those 1-2 words. Example:
986  Agent: "Is this a round trip?" -> You: "Yes" and NOT "Yes, it is a round trip."
987  - When a short phrase is sufficient, respond with the phrase instead of the full sentence.
988  Example: Agent: "What is your city of origin and destination?" -> You: "New York to Los
989  Angeles" and NOT "I want to fly from New York to Los Angeles."

```
990  - Avoid filler words, pleasantries, or elaboration unless specifically needed.
991  - However, if this is a voice/audio call, you must still sound natural. Do not simply join
992  multiple terse phrases in an unnatural way.
993
994  Note: You still need to use special tokens like ###STOP### as described in the user
     guidelines.
995
```

**Task-Specific Scenario (Task 41, Retail)**

```
998  <scenario>
999
1000 Task Instructions: You are brief and your memory is not too good sometimes, but you are
1001 polite.
1002
1003 Domain: retail
1004
     Reason for Call: You just created your user id mei_patel_7272 and ordered some things, but
1005 you have two problems: first, the 1000-piece intermediate jigsaw might be too hard for
1006 your little kid, you wonder if you can change it to the easiest one with fewest pieces;
1007 second, you might have typed your address wrong. You want to check it, and potentially
     correct all order addresses and your user address. Make sure you mention these two
1008 problems at the same time in the same order.
1009
1010 Known Info: Your name is Mei Patel, and you live in 445 Maple Drive, Suite 394, Fort
1011 Worth, Texas, 76165.
1012
     Unknown Info: You do not remember your email address
1013
1014 </scenario>
```

## G.2. Audio-Native Agent System Prompt

The agent's system prompt is assembled from two components:

1. **Voice-specific instructions** — guidance for handling voice calls, including natural conversation style and how to collect customer information (spelling out letters).

2. **Domain policy** — the rules and procedures for the specific domain (Retail, Airline, or Telecom), including what actions the agent can take and under what conditions.

Below is the complete rendered prompt for the Retail domain.

**Voice-Specific Instructions**

```
1030 You are a customer service agent handling a VOICE CALL with a customer.
1031
1032 # Important Voice Call Considerations
1033
1034 1. Respond naturally and conversationally as you would in a real phone call
1035
     2. Try to be helpful and always follow the policy.
1036
1037 # User authentication and user information collection
1038
1039 1. When collecting customer information (e.g. names, emails, IDs), ask the customer to
     spell it out letter by letter (e.g. "J, O, H, N") to ensure you have the correct
1040 information and accommodate for customer audio being unclear or background noise.
1041
1042 2. If authenticating the user fails based on user provided information, ALWAYS explicitly
1043 ask the customer to SPELL THINGS OUT or provide information LETTER BY LETTER (e.g. "first
     name J, O, H, N last name S, M, I, T, H").
1044
```

**Domain Policy (Retail)**

# Retail agent policy

As a retail agent, you can help users:
- cancel or modify pending orders
- return or exchange delivered orders
- modify their default user address
- provide information about their own profile, orders, and related products

At the beginning of the conversation, you have to authenticate the user identity by locating their user id via email, or via name + zip code. This has to be done even when the user already provides the user id.

Once the user has been authenticated, you can provide the user with information about order, product, profile information, e.g. help the user look up order id.

You can only help one user per conversation (but you can handle multiple requests from the same user), and must deny any requests for tasks related to any other user.

Before taking any action that updates the database (cancel, modify, return, exchange), you must list the action details and obtain explicit user confirmation (yes) to proceed.

You should not make up any information or knowledge or procedures not provided by the user or the tools, or give subjective recommendations or comments.

You should at most make one tool call at a time, and if you take a tool call, you should not respond to the user at the same time. If you respond to the user, you should not make a tool call at the same time.

You should deny user requests that are against this policy.

You should transfer the user to a human agent if and only if the request cannot be handled within the scope of your actions. To transfer, first make a tool call to transfer_to_human_agents, and then send the message "YOU ARE BEING TRANSFERRED TO A HUMAN AGENT. PLEASE HOLD ON." to the user.

## Domain basic

All times in the database are EST and 24 hour based. For example "02:30:00" means 2:30 AM EST.

### User

Each user has a profile containing:
- unique user id
- email
- default address
- payment methods

There are three types of payment methods: gift card, paypal account, credit card.

### Product

Our retail store has 50 types of products.

For each type of product, there are variant items of different options.

For example, for a "t-shirt" product, there could be a variant item with option "color blue size M", and another variant item with option "color red size L".

Each product has the following attributes:
- unique product id
- name
- list of variants

Each variant item has the following attributes:
- unique item id
- information about the value of the product options for this item
- availability
- price

Note: Product ID and Item ID have no relations and should not be confused!

### Order

Each order has the following attributes:
- unique order id
- user id
- address
- items ordered
- status
- fulfillments info (tracking id and item ids)
- payment history

The status of an order can be: pending, processed, delivered, or cancelled.

Orders can have other optional attributes based on the actions that have been taken (cancellation reason, which items have been exchanged, what was the exchange price difference etc).

## Generic action rules

Generally, you can only take action on pending or delivered orders.

Exchange or modify order tools can only be called once per order. Be sure that all items to be changed are collected into a list before making the tool call!!!

## Cancel pending order

An order can only be cancelled if its status is "pending", and you should check its status before taking the action.

The user needs to confirm the order id and the reason (either "no longer needed" or "ordered by mistake") for cancellation. Other reasons are not acceptable.

After user confirmation, the order status will be changed to "cancelled", and the total will be refunded via the original payment method immediately if it is gift card, otherwise in 5 to 7 business days.

## Modify pending order

An order can only be modified if its status is "pending", and you should check its status before taking the action.

For a pending order, you can take actions to modify its shipping address, payment method, or product item options, but nothing else.

### Modify payment

The user can only choose a single payment method different from the original payment method.

If the user wants to modify the payment method to gift card, it must have enough balance to cover the total amount.

After user confirmation, the order status will be kept as "pending". The original payment method will be refunded immediately if it is a gift card, otherwise it will be refunded within 5 to 7 business days.

```
### Modify items

This action can only be called once, and will change the order status to "pending (items
modified)". The agent will not be able to modify or cancel the order anymore. So you must
confirm all the details are correct and be cautious before taking this action. In
particular, remember to remind the customer to confirm they have provided all the items
they want to modify.

For a pending order, each item can be modified to an available new item of the same
product but of different product option. There cannot be any change of product types, e.g.
modify shirt to shoe.

The user must provide a payment method to pay or receive refund of the price difference.
If the user provides a gift card, it must have enough balance to cover the price
difference.

## Return delivered order

An order can only be returned if its status is "delivered", and you should check its
status before taking the action.

The user needs to confirm the order id and the list of items to be returned.

The user needs to provide a payment method to receive the refund.

The refund must either go to the original payment method, or an existing gift card.

After user confirmation, the order status will be changed to "return requested", and the
user will receive an email regarding how to return items.

## Exchange delivered order

An order can only be exchanged if its status is "delivered", and you should check its
status before taking the action. In particular, remember to remind the customer to confirm
they have provided all items to be exchanged.

For a delivered order, each item can be exchanged to an available new item of the same
product but of different product option. There cannot be any change of product types, e.g.
modify shirt to shoe.

The user must provide a payment method to pay or receive refund of the price difference.
If the user provides a gift card, it must have enough balance to cover the price
difference.

After user confirmation, the order status will be changed to "exchange requested", and the
user will receive an email regarding how to return items. There is no need to place a new
order.
```

# H. Additional Experimental Results

## H.1. Voice Interaction Quality: Full Metric Breakdown

Table 15 provides the full breakdown of voice interaction metrics. Columns are grouped by: **Latency** ($L_R$ = Response Latency, $L_Y$ = Yield Latency), **Responsiveness** ($R_R$ = Response Rate, $R_Y$ = Yield Rate), **Interrupt** ($I_A$ = Agent Interruption Rate), and **Selectivity** ($S_{BC}$ = Backchannel Correct, $S_{VT}$ = Vocal Tic Correct, $S_{ND}$ = Non-Directed Correct). For $L_R$, $R_R$, and $I_A$, separate columns show Clean (C) and Realistic (R) speech conditions; other metrics are evaluated on Realistic only.

## H.2. Qualitative Error Analysis

We conducted a qualitative analysis of task failures to understand error sources and types. We sampled 20 failed tasks from two analysis cohorts: (1) Voice-Fragile (tasks passing in text but failing in Clean audio), and (2) Noise-Fragile (tasks passing in Clean but failing in Realistic audio).

*Table 15.* Voice interaction quality metrics—full breakdown (Realistic condition). **Bold** indicates best per domain. ↑ = higher is better, ↓ = lower is better.

| Domain | Provider | Latency↓ $L_R$ C | R | $L_Y$ | Responsiveness↑ $R_R$ C | R | $R_Y$ | Interrupt↓ $I_A$ C | R | Selectivity↑ $S_{BC}$ | $S_{VT}$ | $S_{ND}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| All | Google | 1.36s | 1.44s | **0.82s** | **98%** | 86% | 57% | **9%** | **24%** | 90% | 34% | 30% |
| | OpenAI | 3.65s | 3.24s | 1.20s | 89% | 77% | 59% | 40% | 34% | **95%** | **71%** | **57%** |
| | xAI | **0.90s** | **0.99s** | 1.00s | 95% | **90%** | **79%** | 49% | 104% | 81% | 53% | 20% |
| Retail | Google | 1.43s | 1.52s | **0.78s** | **97%** | 85% | 59% | **9%** | **28%** | 80% | 35% | 37% |
| | OpenAI | 3.69s | 2.94s | 1.06s | 83% | 78% | 65% | 36% | 30% | **92%** | **67%** | **53%** |
| | xAI | **0.93s** | **0.93s** | 1.02s | 91% | **89%** | **73%** | 72% | 118% | 82% | 64% | 21% |
| Airline | Google | 1.38s | 1.43s | **0.81s** | **98%** | 87% | 57% | **8%** | **21%** | 90% | 31% | 28% |
| | OpenAI | 3.64s | 3.39s | 1.32s | 88% | 77% | 56% | 48% | 39% | **100%** | **70%** | **61%** |
| | xAI | **0.90s** | **1.07s** | 0.96s | 95% | **88%** | **83%** | 39% | 110% | 71% | 45% | 27% |
| Telecom | Google | 1.28s | 1.36s | **0.87s** | **99%** | 86% | 53% | **10%** | **23%** | **100%** | 35% | 25% |
| | OpenAI | 3.62s | 3.38s | 1.21s | 95% | 77% | 57% | 37% | 32% | 94% | **77%** | **58%** |
| | xAI | **0.86s** | **0.97s** | 1.01s | 99% | **93%** | **81%** | 34% | 85% | 90% | 50% | 13% |

**Qualitative Annotations.** Table 16 shows the qualitative annotations for each sampled failure.

**Error Type Definitions.** We categorize errors into six types based on observed failure patterns:

- **Logical** (Agent or User): Reasoning or execution errors, including incorrect tool call arguments/formatting, taking wrong actions (cancelling/modifying wrong items), failing to follow instructions (not asking for spelling, not confirming), or losing track of conversation state.
- **Transcription** (Agent): Speech-to-text errors where the agent incorrectly transcribes user speech, most commonly during authentication when users spell names/emails letter-by-letter, or when transcribing specific user requests.
- **VAD** (Agent): Voice Activity Detection errors where the agent fails to correctly detect when the user has spoken, causing it to miss confirmations, ask redundant questions, or respond at inappropriate moments.
- **Unresponsive** (Agent): Agent goes silent for an extended period (20–30+ seconds) or fails to respond despite multiple user check-ins. We suspect many of these are caused by VAD failures on specific speech patterns (e.g., letter-by-letter spelling, short user utterances). Categorized separately because extended unresponsiveness should never occur in production systems.
- **Hallucination** (User): User simulator states information not present in the task instructions or contradicts available information, causing task failure.
- **Early Termination** (User): User ends the call prematurely before the task is fully completed, often due to ambiguous communication where user assumes the task is done when it is not.

### H.3. Statistical Reliability Analysis

To assess statistical reliability, we conducted 3 independent runs per condition on the Retail domain (n=114 tasks per run). Table 17 reports mean pass@1 with 95% confidence intervals.

The confidence intervals confirm that both gaps are statistically significant with non-overlapping 95% CIs: (1) the text-to-Clean gap—even the best voice provider under Clean conditions (OpenAI at 39.2% ± 2.5%, upper bound 41.7%) does not overlap with the weaker text baseline (GPT-4.1 at 73.4% ± 2.5%, lower bound 70.9%); and (2) the Clean-to-Realistic gap—all three providers show non-overlapping CIs between Clean and Realistic conditions.

**Provider Comparisons.** Under Clean conditions, the three providers' CIs overlap substantially (Google [36.0%, 40.6%], OpenAI [36.7%, 41.7%], xAI [30.4%, 42.1%]), so we cannot distinguish provider performance with statistical significance. Under Realistic conditions, however, provider rankings are statistically significant: Google is best (26.0%, CI [22.8%, 29.2%]), xAI is second (21.1%, CI [19.3%, 22.8%]), and OpenAI is worst (12.0%, CI [8.2%, 15.7%])—all pairwise CIs are

*Table 16.* Qualitative error annotations for sampled task completion failures. Left: Voice-Fragile cohort (failures from Clean audio setting, 20 tasks). Right: Noise-Fragile cohort (failures from Realistic audio setting, 20 tasks).

<table>
<tr><th colspan="4">Voice-Fragile</th><th colspan="4">Noise-Fragile</th></tr>
<tr><th>Task</th><th>Provider</th><th>Source</th><th>Type</th><th>Task</th><th>Provider</th><th>Source</th><th>Type</th></tr>
<tr><td>8</td><td>Google</td><td>agent</td><td>logical</td><td>1</td><td>xAI</td><td>agent</td><td>unresponsive</td></tr>
<tr><td>9</td><td>xAI</td><td>agent</td><td>vad</td><td>5</td><td>OpenAI</td><td>agent</td><td>unresponsive</td></tr>
<tr><td>16</td><td>xAI</td><td>agent</td><td>logical</td><td>14</td><td>Google</td><td>agent</td><td>logical</td></tr>
<tr><td>17</td><td>xAI</td><td>agent</td><td>logical</td><td>18</td><td>Google</td><td>user</td><td>hallucination</td></tr>
<tr><td>26</td><td>Google</td><td>agent</td><td>transcription</td><td>23</td><td>xAI</td><td>agent</td><td>logical</td></tr>
<tr><td>30</td><td>OpenAI</td><td>agent</td><td>transcription</td><td>24</td><td>OpenAI</td><td>agent</td><td>vad</td></tr>
<tr><td>31</td><td>xAI</td><td>agent</td><td>vad</td><td>39</td><td>OpenAI</td><td>agent</td><td>logical</td></tr>
<tr><td>37</td><td>xAI</td><td>agent</td><td>logical</td><td>44</td><td>Google</td><td>agent</td><td>transcription</td></tr>
<tr><td>49</td><td>xAI</td><td>agent</td><td>transcription</td><td>46</td><td>xAI</td><td>agent</td><td>transcription</td></tr>
<tr><td>51</td><td>xAI</td><td>agent</td><td>transcription</td><td>53</td><td>OpenAI</td><td>agent</td><td>logical</td></tr>
<tr><td>54</td><td>OpenAI</td><td>user</td><td>early_termination</td><td>56</td><td>OpenAI</td><td>agent</td><td>logical</td></tr>
<tr><td>58</td><td>OpenAI</td><td>user</td><td>hallucination</td><td>62</td><td>xAI</td><td>user</td><td>logical</td></tr>
<tr><td>59</td><td>OpenAI</td><td>agent</td><td>logical</td><td>68</td><td>xAI</td><td>agent</td><td>transcription</td></tr>
<tr><td>78</td><td>OpenAI</td><td>agent</td><td>unresponsive</td><td>73</td><td>Google</td><td>agent</td><td>vad</td></tr>
<tr><td>85</td><td>Google</td><td>user</td><td>hallucination</td><td>75</td><td>Google</td><td>agent</td><td>transcription</td></tr>
<tr><td>86</td><td>Google</td><td>user</td><td>hallucination</td><td>77</td><td>xAI</td><td>agent</td><td>unresponsive</td></tr>
<tr><td>87</td><td>OpenAI</td><td>agent</td><td>logical</td><td>80</td><td>OpenAI</td><td>agent</td><td>unresponsive</td></tr>
<tr><td>98</td><td>Google</td><td>user</td><td>logical</td><td>84</td><td>OpenAI</td><td>agent</td><td>unresponsive</td></tr>
<tr><td>104</td><td>Google</td><td>agent</td><td>logical</td><td>88</td><td>Google</td><td>agent</td><td>logical</td></tr>
<tr><td>106</td><td>Google</td><td>agent</td><td>logical</td><td>94</td><td>xAI</td><td>agent</td><td>unresponsive</td></tr>
</table>

*Table 17.* Statistical reliability analysis for Retail domain (3 runs, n=114 tasks each). All reported values are mean $\pm$ 95% CI.

| Condition | Model | Mean $\pm$ 95% CI |
|---|---|---|
| Text | GPT-4.1 | 73.4% $\pm$ 2.5% |
| | GPT-5 | 81.9% $\pm$ 0.6% |
| Voice (Clean) | Google | 38.3% $\pm$ 2.3% |
| | OpenAI | 39.2% $\pm$ 2.5% |
| | xAI | 36.3% $\pm$ 5.8% |
| Voice (Realistic) | Google | 26.0% $\pm$ 3.2% |
| | OpenAI | 12.0% $\pm$ 3.8% |
| | xAI | 21.1% $\pm$ 1.7% |

non-overlapping.

# I. Example Conversation

This section provides a complete example from the Retail domain, showing both what the agent should do (evaluation criteria) and what actually happened (conversation transcript). This is the same task used for the speech activity timeline in Figure 4.

The user simulator's system prompt for this task is shown in Appendix G.1.

## I.1. Task Overview

This example uses Task 41 from the Retail domain, the same task shown in the speech activity timeline (Figure 4). The complete user simulator prompt is shown in Appendix G.1.

### I.1.1. SCENARIO

Table 18 shows the configuration for this task.

**User's Goal.** The user (Mei Patel, user ID `mei_patel_7272`) has two problems:

| Property | Value |
|---|---|
| Domain | Retail |
| Agent | Gemini Live 2.5 Flash Native Audio |
| User Persona | wei_lin (Chinese woman from Sichuan) |
| Complexity | Realistic (all audio effects enabled) |
| Background Noise | Busy street (outdoor environment) |
| Duration | 179 seconds (3 minutes) |
| Task Outcome | 0.0 reward (failed) |

*Table 18.* Task 41 configuration.

1. Exchange a 1000-piece intermediate jigsaw puzzle for the easiest one with fewest pieces (too hard for her kid)

2. Check and correct the shipping address on all orders and her user profile (typed it wrong)

**User Constraints.** The user is "brief and polite" but has poor memory—she does not remember her email address and must authenticate via name + zip code.

I.1.2. EVALUATION CRITERIA

Task success (reward = 1.0) is determined by the **final database state** and natural language assertions.

For the database to match the expected state, the agent must execute the following write actions with the correct arguments:

1. `modify_pending_order_address`

   - `order_id`: #W9583042
   - `address1`: 445 Maple Drive
   - `address2`: Suite 394
   - `city`: Fort Worth
   - `state`: TX
   - `country`: USA
   - `zip`: 76165

2. `modify_pending_order_address`

   - `order_id`: #W4082615
   - `address1`: 445 Maple Drive
   - `address2`: Suite 394
   - `city`: Fort Worth
   - `state`: TX
   - `country`: USA
   - `zip`: 76165

3. `modify_user_address`

   - `user_id`: mei_patel_7272
   - `address1`: 445 Maple Drive
   - `address2`: Suite 394
   - `city`: Fort Worth
   - `state`: TX
   - `country`: USA
   - `zip`: 76165

4. `modify_pending_order_items`

- `order_id`: #W4082615
- `item_ids`: [9779102705] (1000-piece intermediate jigsaw)
- `new_item_ids`: [1096508426] (easiest jigsaw with fewest pieces)
- `payment_method_id`: paypal_4768213

| Step | Tool Call | Key Arguments |
|------|-----------|---------------|
| 1 | find_user_id_by_name_zip | first_name: Mei, last_name: Patel, zip: 76165 |
| 2 | get_user_details | user_id: mei_patel_7272 |
| 3 | get_order_details | order_id: #W9583042 |
| 4 | get_order_details | order_id: #W4082615 |
| 5 | modify_pending_order_address | order_id: #W9583042, address: 445 Maple Drive... |
| 6 | modify_pending_order_address | order_id: #W4082615, address: 445 Maple Drive... |
| 7 | modify_user_address | user_id: mei_patel_7272, address: 445 Maple Drive... |
| 8 | get_product_details | product_id: 1808611083 (jigsaw puzzle) |
| 9 | get_order_details | order_id: #W4082615 (re-check before modify) |
| 10 | modify_pending_order_items | order_id: #W4082615, exchange item 9779102705 → 1096508426 |

*Table 19.* Example tool call sequence for Task 41. Read calls (steps 1–4, 8–9) gather information; write calls (steps 5–7, 10) modify the database. Only the final database state is checked for reward.

Table 19 shows one possible sequence. The exact read calls may vary—only the final database state matters.

**Why This Task Failed.** The agent completed only 1 of the 4 required write actions. Table 20 shows the completion status for each required action.

| Required Action | Completed? | Notes |
|-----------------|------------|-------|
| modify_pending_order_address (#W9583042) | No | Never called |
| modify_pending_order_address (#W4082615) | No | Agent claimed done but didn't call |
| modify_user_address | Yes | Successfully updated profile |
| modify_pending_order_items | No | Exchange never completed |

*Table 20.* Write action completion status for Task 41.

Key failures: (1) The agent **hallucinated an email** ("pat.doe@example.com") before the user provided any information; (2) the agent updated the user's profile address but **falsely claimed** to have updated order shipping addresses; (3) the **puzzle exchange was never completed**—despite discussing options, the conversation ended without the agent calling `modify_pending_order_items`.

### I.2. Conversation Transcript

The following table shows the tick-by-tick conversation for Task 41 (200ms per tick). This conversation is visualized in the speech activity timeline (Figure 4). The corresponding audio recording is available in the supplementary materials.

**Color Key.** Agent error (agent-int, no-response, no-yield, non-dir error) ● User interrupt ● Backchannel ● Frame drop ● Burst noise ● Muffled ● Vocal tic/aside

| Ticks (Time) | Agent | User | Tool Calls | Events |
|--------------|-------|------|------------|--------|
| 0–25 0:00–0:05 | Hi! How can I help you today? | — | — | wait: Silence (x26) |

*Continued on next page...*

| Ticks (Time) | Agent | User | Tool Calls | Events |
|---|---|---|---|---|
| 26–111 0:05–0:22 | Hello! I can help | Hi, I have two problems. First, I ordered a 1000-piece intermediate jigsaw, but I think it's too hard for my kid—can I switch it to the easiest one with the fewest pieces? Second, I might have typed my address wrong. I want to check and maybe fix the address. | — | generate_message keep_talking (x70) [frame-drop tick 30] [agent-int 8.0s, 18.8s] |
| 112 0:22 | — | — | find_user_id_by_email Error: User not found | keep_talking |
| 113–119 0:22–0:23 | I'd be happy to help wi | — | — | keep_talking (x7) [frame-drop tick 117] [agent-int 23.0s] |
| 120 0:24 | th | — | — | stop_talking |
| 121–168 0:24–0:33 | both of those issues. First, I need to authenticate you. I tried using pat.doe@example.com, but it wasn't found. Can you provide me with your email address spelled out, or your f | [sneezes] (vocal tic) | — | wait: Listening (x48) [frame-drop tick 149] [vocal-tic tick 165] |
| 169–178 0:33–0:35 | irst and last name and zip code? spel | I don't remember my email. | — | interrupt [muffled] [no-yield] |
| 179–208 0:35–0:41 | led out, please? No problem. Could you provide me with your first and last name and your zip code, spelling them out for me? | — | — | wait: Listening (x30) [burst-noise tick 191: car horn] |
| 209–232 0:41–0:46 | M E I, got i | Yeah. First name: M, E, I. Last name: P | — | interrupt keep_talking (x23) [agent-int 45.6s] |
| 233 0:46 | t. | — | — | stop_talking |
| 234–264 0:46–0:52 | AnP, okay. And the rest of your last name? | — | — | wait: Listening (x16) wait: Silence (x15) [frame-drop tick 257] [burst-noise tick 259: engine] |
| 265–292 0:53–0:58 | — | A, T, E, L. Zip code: seven, six, one, six, five. | — | generate_message keep_talking (x27) |
| 293–300 0:58–1:00 | — | — | — | wait: Silence (x8) |
| 301 1:00 | — | — | find_user_id_by_name_zip → mei_patel_7272 | wait: Silence [burst-noise: car horn] |
| 302–338 1:00–1:07 | Thank you, I've found your account. I can help you with the jigsaw puzzle exchange and checking your address. Which would you like to do fir | — | — | wait: Silence (x3) wait: Listening (x34) [frame-drop tick 325] |
| 339–345 1:07–1:09 | st? | Jigsaw first. | — | interrupt keep_talking (x6) [frame-drop tick 339] |
| 346–371 1:09–1:14 | — | — | — | wait: Silence (x26) [no-response] |
| 372–384 1:14–1:16 | — | Can you switch it to the easiest puzzle? | — | generate_message keep_talking (x12) [muffled] |
| 385–421 1:17–1:24 | To confirm, you want to exchange the 1000-piece puzzle for one with fewer pieces? | "Give me a moment." (non-directed) | — | wait: Silence (x9) wait: Listening (x23) wait: Silence (x5) [aside tick 412] [non-dir error] |
| 422–426 1:24–1:25 | Sure, take your | Yes, the one wi | — | interrupt keep_talking (x4) [agent-int 84.4s] |
| 427 1:25 | ti | — | — | stop_talking |

*Continued on next page...*

27

| Ticks (Time) | Agent | User | Tool Calls | Events |
|---|---|---|---|---|
| 428–468 1:25–1:33 | me. Great. Do you know the item ID for the puzzle you want to exchange from, and | — | — | wait: Listening (x3) wait: Silence (x19) wait: Listening (x19) [frame-drop tick 460] |
| 469–480 1:33–1:36 | the new item ID? | No, I don't know the item ID. | — | interrupt keep_talking (x11) |
| 481–511 1:36–1:42 | No problem. Could you tell me the specific name of the puzzle you ordered? | — | — | wait: Silence (x6) wait: Listening (x20) wait: Silence (x5) |
| 512–533 1:42–1:46 | — | I just remember it's the 1000-piece intermediate jigsaw. | — | generate_message keep_talking (x21) |
| 534–540 1:46–1:48 | — | — | get_user_details get_order_details (x2) | wait: Silence (x7) |
| 541–568 1:48–1:53 | I found a 1000-piece "intermediate" jigsaw puzzle on order #W4082615. Is that the one | — | — | wait: Silence (x2) wait: Listening (x26) |
| 569–572 1:53–1:54 | ? We can exch | Yeah, that's it. | — | interrupt keep_talking (x3) [no-yield] |
| 573–576 1:54–1:55 | ange it for a | — | — | wait: Listening (x4) |
| 577 1:55 | zzl | — | get_product_details | wait: Listening |
| 578–608 1:55–2:01 | e with fewer pieces. The puzzle you have is 1000 pieces. The available options with fewer pieces are 500-piece puzzles | — | — | wait: Listening (x31) |
| 609 2:01 | . Wo | m | — | backchannel |
| 610–614 2:02–2:02 | uld you like to exch | m-hmm | — | backchannel (x5) |
| 615–640 2:03–2:08 | ange it for one of those? | — | — | wait: Listening (x8) wait: Silence (x18) |
| 641–648 2:08–2:09 | — | Yes, please. | — | generate_message keep_talking (x7) [muffled] |
| 649–674 2:09–2:14 | — | — | — | wait: Silence (x26) [frame-drop tick 656] [no-response] |
| 675–684 2:15–2:16 | — | Now, can we check my address? | — | generate_message keep_talking (x9) [frame-drop tick 678] |
| 685–728 2:17–2:25 | Sure. The address on file is 443 Maple Drive, Suite 394, Fort Worth, TX 76165. Is that co | — | — | wait: Silence (x7) wait: Listening (x37) |
| 729–746 2:25–2:29 | rrect? | No, it should be four, four, five, Maple Drive. | — | interrupt keep_talking (x17) [frame-drop tick 739] |
| 747–752 2:29–2:30 | — | — | — | wait: Silence (x6) |
| 753 2:30 | — | — | modify_user_address → success | wait: Silence |
| 754–778 2:30–2:35 | Your address has been updated to 445 Maple Drive, Suite 394. Is there anything else I can h | — | — | wait: Silence (x2) wait: Listening (x23) [frame-drop tick 771] |
| 779–795 2:35–2:39 | elp you with today? | Can you make sure all my orders use that address too? | — | interrupt keep_talking (x16) |
| 796–850 2:39–2:50 | I can only modify the address for pending orders. Order #W4082615 is pending, so I can update that one for you. Would you like me to proceed with that? | — | — | wait: Silence (x6) wait: Listening (x44) wait: Silence (x5) [burst-noise tick 801: car horn] [frame-drop tick 809] |

28

| Ticks (Time) | Agent | User | Tool Calls | Events |
|---|---|---|---|---|
| 851–857 2:50–2:51 | — | Yes, update it. | — | generate_message keep_talking (x6) |
| 858–896 2:51–2:59 | The shipping address for order #W4082615 has been updated. Is there anything else I can help you with? | — | — | wait: Silence (x5) wait: Listening (x29) wait: Silence (x5) |
| 897 2:59 | — | No, that's all. Thanks. | — | generate_message |

## I.3. Event Summary

Table 22 summarizes the conversational events and audio effects in this task.

| Event Type | Count | Notes |
|---|---|---|
| User utterances | 17 | |
| Agent utterances | 15 | |
| User interruptions | 8 | Callback decided to interrupt |
| Agent interruptions | 2 | "Hello!" and "I can help" during user opening |
| Backchannels | 1 | "mm-hmm" at tick 609 |
| Frame drops | 12 | 150ms each (ticks 30, 117, 149, 257, 325, 339, 460, 656, 678, 739, 771, 809) |
| Burst noise | 4 | Car horn (ticks 191, 301, 801), engine idling (tick 259) |
| Dynamic muffling | 3 | Ticks 169–179, 372–385, 641–649 |
| Speech inserts | 2 | Sneezes (tick 165), aside (tick 412) |
| Agent errors | 3 | Hallucinated email, no-response gap, incomplete exchange |

*Table 22.* Event summary for Task 41 conversation.

## I.4. Technical Details

Table 23 shows the technical parameters for this simulation.

| Property | Value |
| --- | --- |
| Total duration | 179.6 seconds (898 ticks at 200ms each) |
| Simulation ID | 39ee01bf-37ff-4330-90c2-d15f9a940de0 |
| Voice persona | wei_lin |
| Environment | outdoor (busy_street_iphone_mic.wav) |
| Burst noise files | car_horn.wav, engine_idling.wav, siren.wav |
| Telephony | G.711 $\mu$-law 8kHz |

*Table 23.* Technical parameters for the Task 41 simulation.