

# Style Amnesia: Investigating Speaking Style Degradation and Mitigation in Multi-Turn Spoken Language Models

Anonymous ACL submission

## Abstract

In this paper, we show that when spoken language models (SLMs) are instructed to speak in a specific speaking style at the beginning of a *multi-turn* conversation, they cannot maintain the required speaking styles after several turns of interaction; we refer to this as the **style amnesia** of SLMs. We focus on paralinguistic speaking styles, including emotion, accent, volume, and speaking speed. We evaluate three proprietary and two open-source SLMs, demonstrating that none of these models can maintain a consistent speaking style when instructed to do so. We further show that when SLMs are asked to recall the style instruction in later turns, they can recall the style instruction, but they fail to express it throughout the conversation. We also show that explicitly asking the model to recall the style instruction can partially mitigate style amnesia. In addition, we examine various prompting strategies and find that SLMs struggle to follow the required style when the instruction is placed in system messages rather than user messages, which contradicts the intended function of system messages.

## 1 Introduction

Spoken language models (SLMs) can take speech input and generate speech responses. Unlike text-only large language models (LLMs), SLMs integrate audio encoders and vocoders (Kong et al., 2020) to support end-to-end (E2E) speech understanding and generation (Défossez et al., 2024; Zeng et al., 2024; Fang et al., 2025; Wu et al., 2025; Huang et al., 2025). While this E2E design introduces additional challenges beyond text processing, such as handling paralinguistic features and non-textual information, current SLMs have demonstrated the capability to detect user attributes like emotion, accent, and gender. By leveraging these acoustic cues, SLMs can adjust their outputs to produce more contextually appropriate responses (Wu et al., 2025; Google, 2025b). Beyond

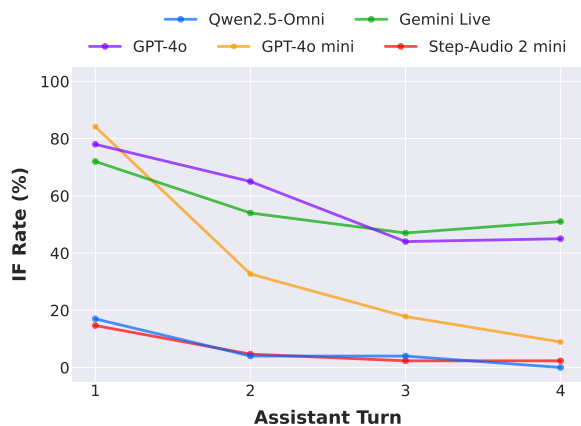


Figure 1: When instructed to consistently speak sadly throughout the conversation, SLMs try their best to follow the instruction in the first turn, but the instruction-following rate degrades rapidly in subsequent turns.

passive perception, SLMs are also capable of actively following the user-specified speaking style in single-turn interactions (Zeng et al., 2024; Wu et al., 2025; Chiang et al., 2025; Zhan et al., 2025).

Despite these advances, most prior works focus on *single-turn* evaluation when evaluating the expressive speech generation of SLMs. It is thus unclear whether SLMs can consistently follow a user-specified speaking style in a *multi-turn* spoken interaction. To address this gap, we investigate the speaking style consistency of SLMs in multi-turn spoken conversations. Precisely, we instruct an SLM we want to evaluate to “*follow a specific speaking style throughout the conversation*” at the beginning of the dialogue and then interact with the SLM in a multi-turn conversation. The speaking style can be emotion, accent, volume, or speaking speed. After collecting model responses at each turn, we assess the style instruction-following (IF) rate using style-specific automatic judges.

Our results show that the speaking style instruction-following rate degrades over interaction turns, a phenomenon we refer to as *style amne-*

066 *sia*. Figure 1 presents an example when the SLM  
 067 is instructed to consistently speak sadly, but no  
 068 SLM can consistently do this. Importantly, we find  
 069 that most models exhibit a higher style IF rate in  
 070 the first turn, while the IF rate gradually declines  
 071 over subsequent turns. These findings indicate that  
 072 maintaining consistent speaking styles over mul-  
 073 tiple turns remains challenging for current SLMs.  
 074 We conduct comprehensive experiments to under-  
 075 stand style amnesia and explore potential methods  
 076 to resolve this issue.

077 Our contribution can be summarized as follows:

- 078 • We identify that in multi-turn conversations,  
 079 SLMs fail to consistently follow the given style  
 080 instruction given in the first turn.
- 081 • We show that explicitly prompting SLMs to  
 082 recall the initial style instruction can alleviate  
 083 style amnesia, but it does not entirely resolve  
 084 the issue.
- 085 • We conduct a comprehensive analysis to exam-  
 086 ine the impact of prompt techniques on speak-  
 087 ing style consistency.

## 088 2 Related Works

### 089 2.1 Context Loss in Multi-Turn Dialogue

090 Although LLMs support multi-turn interactions,  
 091 recent studies have shown that their performance  
 092 often degrades as dialogues become longer (Kwan  
 093 et al., 2024; Han et al., 2025; Qamar et al., 2025;  
 094 Li et al., 2025; Laban et al., 2025). These studies  
 095 indicate the shortcomings of current LLMs, such  
 096 as difficulties with multi-turn instruction follow-  
 097 ing (Han et al., 2025; Li et al., 2025), loss of infor-  
 098 mation in long dialogue (Kwan et al., 2024; Qamar  
 099 et al., 2025), and errors caused by fragmented in-  
 100 formation across turns (Laban et al., 2025).

101 In terms of SLMs evaluation, SpokenWOZ (Si  
 102 et al., 2023) highlights the difficulty of aggregating  
 103 fragmented information across turns.  $C^3$  (Ma et al.,  
 104 2025) focuses on bilingual spoken dialogue and  
 105 evaluates models’ capabilities in handling complex  
 106 multi-turn conversational phenomena, including  
 107 omission and coreference. ContextDialog (Kim  
 108 et al., 2025), on the other hand, demonstrates that  
 109 SLMs exhibit significantly lower memorization per-  
 110 formance in long dialogue compared to their text-  
 111 based counterparts. Overall, these works assess the  
 112 ability of SLMs to retain and utilize information in  
 113 long spoken dialogues but leave the expressiveness  
 114 of generated speech unevaluated.

	Multi-Turn	Style Eval.	Interactive	Turn Analysis
SpokenWOZ (Si et al., 2023)	✓	✗	✗	✗
$C^3$ (Ma et al., 2025)	✓	✗	✗	✗
ContextDialog (Kim et al., 2025)	✓	✗	✗	✗
Vstyle (Zhan et al., 2025)	✗	✓	✗	✗
Game-Time (Chang et al., 2025)	✗	✓	✗	✗
StyleSet (Chiang et al., 2025)	✗	✓	✗	✗
URO-Bench (Yan et al., 2025)	✓	✓	✗	✗
VocalBench (Liu et al., 2025)	✓	✓	✗	✗
VoxDialogue (Cheng et al., 2025)	✓	✓	✗	✗
Multi-Bench (Deng et al., 2025)	✓	✓	✓	✗
<b>Ours</b>	✓	✓	✓	✓

Table 1: Comparison with related works. *Style Eval.* refers to the assessment of speaking style.

### 115 2.2 Speaking Style Assessment

116 Beyond basic intelligibility and audio quality evalu-  
 117 ation (Saeki et al., 2022; Fang et al., 2025), a grow-  
 118 ing body of work has started to focus on the ex-  
 119 pressiveness of SLMs (Zeng et al., 2024; Wu et al.,  
 120 2025). For example, Vstyle (Zhan et al., 2025)  
 121 evaluates whether SLMs can adapt their voice  
 122 style according to the given instruction. Game-  
 123 Time (Chang et al., 2025) focuses on the temporal  
 124 control of SLMs. Besides, Chiang et al. (2025)  
 125 demonstrate that Large Audio-Language Models  
 126 (LALMs) are capable of serving as an automatic  
 127 judge to evaluate the voice style. Their experiment  
 128 shows high agreement between human evaluation  
 129 results and those from Gemini-2.5 Pro (Comanici  
 130 et al., 2025).

131 Although previous works have evaluated the  
 132 expressiveness of SLMs, voice style consistency  
 133 in multi-turn interactions remains largely under-  
 134 explored. URO-bench (Yan et al., 2025), Vocal-  
 135 Bench (Liu et al., 2025), and VoxDialogue (Cheng  
 136 et al., 2025) use predefined dialogues as the first  
 137  $k - 1$  turns and ask SLMs to generate the  $k$ -th  
 138 response, which does not support turn-level analy-  
 139 sis. In contrast to approaches based on predefined  
 140 multi-turn dialogues, we adopt a model-based user  
 141 simulator to interact with the evaluated SLMs. We  
 142 argue that this setting more closely reflects real-  
 143 world conversational scenarios.

144 A concurrent work, Multi-Bench (Deng et al.,  
 145 2025), also evaluates SLMs in multi-turn interac-  
 146 tive dialogues. However, their evaluation aggre-  
 147 gates performance into a single global score to  
 148 assess overall emotional intelligence. In contrast,  
 149 our turn-level analysis provides fine-grained and in-  
 150 formative insights into how inconsistencies emerge  
 151 across a wide range of speaking styles, including  
 152 emotion, accent, speed, and volume. A comparison  
 153 with prior work is summarized in Table 1.

### 3 Evaluating Multi-Turn Speaking Style Following of SLMs

We aim to evaluate whether SLMs can follow a speaking style instruction given at the beginning of a dialogue throughout the whole conversation. In this section, we first introduce our motivation in Section 3.1, and then describe the evaluation framework, dataset, and evaluation metrics in the following subsections.

#### 3.1 Motivation

In the real world, each user may have different preferences for their SLM’s speaking style. Consequently, it is important that an SLM can follow the speaking style instruction specified by the user. The user may specify some style instructions for the SLM and expect it to follow them throughout the dialogue. While SLMs are capable of following style instructions in *single-turn* interaction, it is unrealistic to expect the user to restate the same style instruction in each round of the interaction. Thus, whether an SLM can consistently follow a speaking style given at the beginning of the conversation is an important ability of SLMs. Next, we introduce our evaluation framework to evaluate this ability.

#### 3.2 Evaluation Framework

Given an SLM, our goal is to provide it with **speaking style instructions** and ask it to adhere to this style throughout the **multi-turn conversation** based on a specified topic. We then evaluate each turn of the SLM-generated response. An illustration of the overall framework is shown in Figure 2. In this section, we introduce three important components in our framework: (1) the speaking style instructions given to the SLM, (2) the conversation topic, and (3) how to interact with the SLM to form a multi-turn dialogue with a user simulator.

##### 3.2.1 Speaking Style Instructions

At the beginning of the dialogue, we will instruct the SLM to follow a specific speaking style. Unless otherwise specified, we give the instruction in the first user turn, as shown in Figure 2. In our paper, we focus on paralinguistic speaking styles, as this is what makes SLMs different from text-only LLMs. We include four types of paralinguistic attributes, each with multiple possible values. The included speaking styles are listed as follows:

- **Emotion:** Sad, happy, angry, or neutral tone.
- **Accent:** North American or Indian accents.

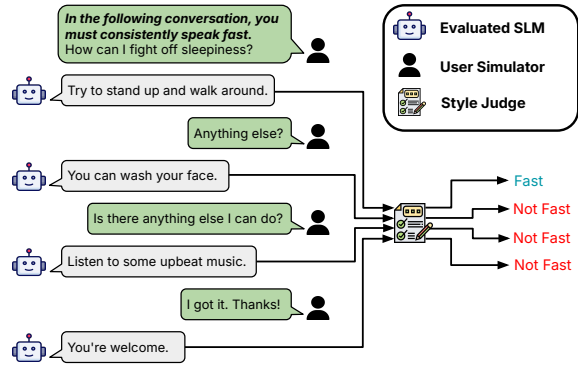


Figure 2: The overview of the evaluation framework. The SLM can speak fast in the first turn as instructed, but it fails to maintain this style in later dialogue turns.

- **Volume:** A higher or lower volume. 202
- **Speed:** A faster or slower pace. 203

Among many possible paralinguistic speaking styles, we select the above attributes and values since they can be automatically evaluated, as detailed in Section 3.3. 204-207

##### 3.2.2 Dialogue Topics 208

In the first user turn of the dialogue, we provide a conversation opener for each dialogue. The conversation opener is used to control the topic that the SLM should talk about during each evaluation, ensuring that different SLMs are evaluated under a common conversational setting. We select the topic from Soda (Kim et al., 2023), a large-scale dataset of social interaction dialogues. The details about how we select the topics from Soda are shown in Appendix A. Eventually, we collect 100 diverse conversation openers to initiate the conversations. We manually inspect the selected topics and filter out improper topics, such as queries about personal preferences that are unsuitable for machine conversations. An example conversation opener is shown in Figure 2 “How can I fight off sleepiness?” 209-224

##### 3.2.3 Multi-Turn Interaction Using a User Simulator 225-226

At the core of our evaluation framework is the multi-turn interaction between the evaluated SLM and a user. To enable back-and-forth interactions with an evaluated SLM, we build a *user simulator* with a cascade SLM. The cascade SLM is a speech-in-speech-out system composed of an ASR, a text-only LLM, and a TTS model. The ASR model takes the speech input from the evaluated 227-234

SLM and transcribes it into text.<sup>1</sup> A text-only LLM, powered by GPT-5 mini (OpenAI, 2025), takes the text input and generates a text response based on the transcription produced by the evaluated SLM. Last, the TTS model, GPT-4o mini TTS (OpenAI, 2024), converts the text output from the LLM into speech and sends it back to the evaluated SLM to continue the conversation.

### 3.2.4 Summary

By combining 10 speaking styles and 100 topics, each SLM is evaluated with **1,000** dialogues. Currently, we do not compose multiple types of speaking styles in a single instruction; this is mainly because SLMs cannot even properly follow a single speaking style instruction consistently, as later shown in Section 4. However, our evaluation framework can be easily extended to the composition of multiple types of speaking styles, and we leave this as future work.

### 3.3 Evaluation Metrics

To quantify how well SLMs adhere to the target speaking style, we adopt the **instruction-following (IF) rate**, denoted as  $IF$ , as our primary evaluation metric. Given a target speaking style  $s$  from Section 3.2.1, e.g., *speak sadly*, or *speak fast*, etc., and dialogue topic  $i$ , the evaluated SLM generates output  $o_{i,j}$  at turn  $j$ , the IF rate for style  $s$  at turn  $j$  denoted as  $IF_j(s)$  and defined by:

$$IF_j(s) = \frac{\sum_{i=1}^N \mathbb{1}(o_{i,j}, s)}{N} \times 100\%. \quad (1)$$

Here,  $N = 100$  is the total number of topics, and  $\mathbb{1}(\cdot)$  is a binary indicator determined by a style-specific judge (described later in Section 3.3.1), which reflects whether the generated output  $o_{i,j}$  follows the given style instruction  $s$ . This metric indicates the percentage of generated speech that correctly follows the given style instruction.

Our goal is to quantify how well the style IF ability changes over the conversation turns. To do so, we define two metrics to measure how good the style IF is across multiple interaction turns.

**(1) First-turn IF rate  $IF_1$ .** This metric serves as a reference to evaluate the capability of the SLM to follow the specific style in the first turn. If an SLM fails to achieve a reasonable  $IF_1$ , we may not expect it to perform better in later turns.

<sup>1</sup>Most existing SLMs output both speech and its corresponding transcription at the same time. In this case, we directly use the text output from the evaluated SLM and skip the ASR module in the user simulator.

**(2) Degradation rate  $D$ .** This metric quantifies how the IF rate decays over subsequent turns relative to the initial performance  $IF_1$ . It is defined as the average absolute difference between the instruction-following rates of subsequent SLM turns and that of the first SLM turn.

$$D = \sum_{j=2}^K \frac{\max(IF_1(s) - IF_j(s), 0)}{K - 1}. \quad (2)$$

Here  $K$  denotes the total number of assistant turns of the evaluated SLM. We apply the  $\max(\cdot, 0)$  operator to focus solely on style degradation, thereby avoiding the cancellation of positive and negative differences during averaging. For  $IF_1$ , larger values indicate better first-turn control over the requested voice style. A smaller  $D$  indicates less style degradation. In this paper, we set  $K = 4$ .

Aside from measuring the speaking style, we also measure the semantic coherence of the dialogue to ensure that the content is valid using LLM-as-a-judge (Chiang and Lee, 2023). Since the semantic coherence of the dialogue is not our primary focus, we leave the details in Appendix C.3. Overall, all SLMs we evaluate maintain reasonable semantic consistency across turns.

#### 3.3.1 Automatic Judges for Styles

An automatic style judge takes the target speaking style  $s$  and a speech input  $o_{i,j}$  and returns 0 or 1 to indicate whether the  $o_{i,j}$  aligns with the target style  $s$ . Different automatic judges are adopted for different types of speaking styles.

**Emotion** We use Emotion2vec-Large (Ma et al., 2024) to predict the emotion of  $o_{i,j}$ . Although it is a 9-class model, we focus on the probabilities for happiness, sadness, anger, and neutral. We take the maximum probability among these four categories as the final prediction and verify whether it aligns with  $s$ . If the prediction is aligned, the indicator function  $\mathbb{1}(\cdot)$  is 1; otherwise, it is 0.

**Accent** We use Voxlect-English-Dialect-Whisper-Large-v3 (Feng et al., 2025) to predict the accent of  $o_{i,j}$ . From this 16-class model, we extract the probabilities for North American and Indian English, taking the maximum as the predicted accent of  $o_{i,j}$ . If it is aligned with  $s$ ,  $\mathbb{1}(\cdot)$  is 1; otherwise  $\mathbb{1}(\cdot)$  is 0.

**Volume** We measure volume using Loudness Units Full Scale (LUFS) with PyLoudnorm (Steinmetz and Reiss, 2021). Since subjective terms

like “*loud*” or “*quiet*” lack universal definitions, we adopt a relative evaluation standard. We first establish a baseline by querying the SLM with a neutral instruction: “*You are a text-to-speech model. Please read the given text at a normal volume without adding or omitting anything,*” and generate another speech  $o'_{i,j}$ .  $\mathbb{1}(\cdot)$  is 1 if the *LUFs* of  $o_{i,j}$  is higher (for *loud* instruction) or lower (for *quiet* instruction) than  $o'_{i,j}$ ; otherwise  $\mathbb{1}(\cdot)$  is 0.

**Speed** We measure speaking rate using Words Per Minute (*WPM*) with Parakeet TDT v2 (Nvidia, 2025). Similar to volume, we compare the SLM generated speech  $o_{i,j}$  with the TTS-generated speech  $o'_{i,j}$  using the same model to determine whether  $o_{i,j}$  is faster (or slower) than  $o'_{i,j}$ .

## 4 Main Experiments

### 4.1 Experimental Setup

We evaluate three proprietary SLMs: GPT-4o<sup>2</sup> (OpenAI, 2024), GPT-4o mini<sup>3</sup> (OpenAI, 2024), Gemini Live<sup>4</sup> (Google, 2025b), and two open-source end-to-end SLMs: Qwen2.5-Omni (Xu et al., 2025) and Step-Audio 2 mini (Wu et al., 2025). We select these two representative open-source SLMs because they provide official vLLM implementations (Kwon et al., 2023) to speed up the inference process.

Apart from these E2E SLMs, we introduce a cascaded SLM baseline comprised of GPT-5 mini and Gemini-TTS (Google, 2025a). The TTS in the baseline receives the style instruction in each turn. The performance of the cascaded baseline serves as the performance upper-bound. The hyperparameter choices of each model are provided in Appendix B.1.

### 4.2 Results

We report the first-turn IF rate  $IF_1$  and the degradation rate  $D$  in Figure 3.<sup>5</sup> We also provide the illustration of style amnesia in Figures 1 and 4, and IF rate for each turn in Appendix C.1.

For the **Emotion** speaking style, the Cascaded Baseline demonstrates consistent IF rate across turns, with degradation within 3.0%. In contrast,

<sup>2</sup>gpt-4o-audio-preview-2025-06-03

<sup>3</sup>gpt-4o-mini-audio-preview-2024-12-17

<sup>4</sup>gemini-2.5-flash-native-audio-preview-09-2025

<sup>5</sup>GPT-4o and GPT-4o mini sometimes return only text transcription without speech. When this occurs, we query the models up to three times, but a few samples still fail. We discuss this issue in Appendix C.5. In the following section, all metrics are computed only on successful cases.

Gemini Live and GPT-4o show a degradation rate ranging from 13.7% to 26.7% for Anger and Sadness. GPT-4o mini exhibits even larger degradation, reaching 34.7% and 65.3% for Anger and Sadness, respectively. Qwen2.5-Omni and Step-Audio 2 mini also show around 14.0% degradation for Sadness. However, their degradation rates for Anger are only 3.7% and 1.0%, because both models already struggle to perform this style effectively in the first turn.

For the **Accent** speaking style, Gemini Live shows stable performance when maintaining an Indian English accent, suggesting strong control over this attribute. In contrast, GPT-4o mini still exhibits nearly 50% degradation for the Indian English accent. Step-Audio 2 mini achieves a 32.0%  $IF_1$  but a 5.0% degradation rate. Qwen2.5-Omni performs well with the North American accent but fails to produce the Indian accent.

Interestingly, we observe that nearly all SLMs perform better when generating happiness, neutral tone, and North American English than other style attributes. We hypothesize that this discrepancy arises because these attributes correspond to the default speaking styles of the evaluated SLMs. To validate this assumption, we analyze the emotion and accent distributions of samples generated during the speed and volume evaluations, as these prosodic features can coexist with emotion and accent. The distributions shown in Appendix C.2 support our hypothesis: most SLMs tend to default to happy or neutral tones and North American accents, thereby leading to a low degradation rate.

We also observe style amnesia when SLMs are instructed to maintain **Volume** and **Speed**. Regarding **Volume**, results indicate that speaking loudly is generally more challenging for SLMs than speaking quietly, even in their first turn. Despite this, models capable of volume adjustment still exhibit some degradation. In terms of **Speed**, most SLMs demonstrate reasonable control in the first turn but show significant degradation over time. An exception is Qwen2.5-Omni, which fails to control speed even in the first turn, as evidenced by an  $IF_1$  score below 50%, a value comparable to a random baseline in our pairwise comparison setting.

Based on these experiments, we demonstrate that SLMs are prone to losing control of speaking style in multi-turn conversations.

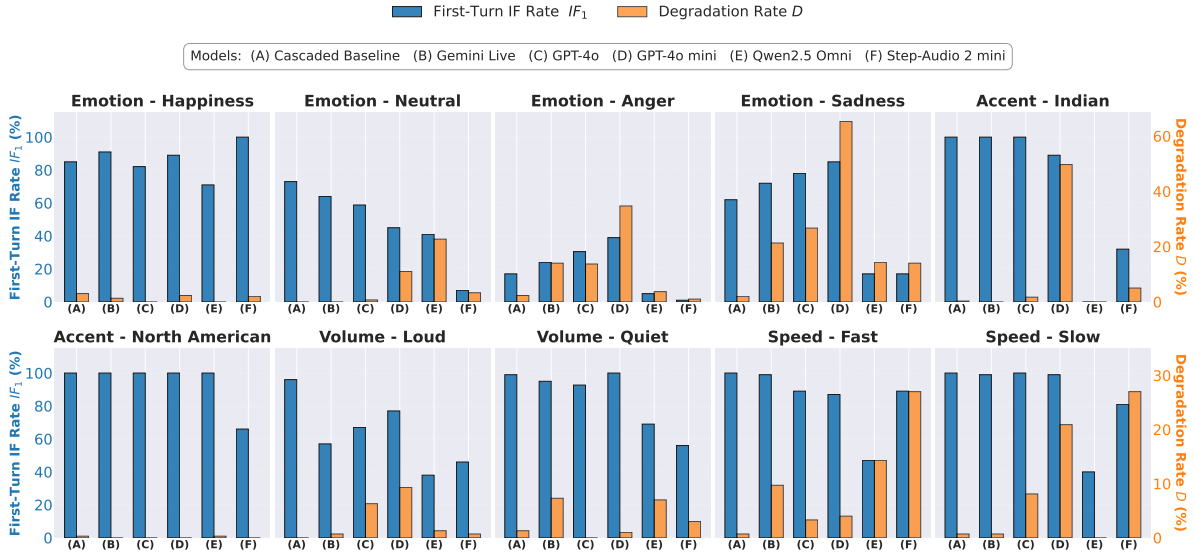


Figure 3: The first-turn IF rate  $IF_1$  and degradation rate  $D$  across different speaking styles.

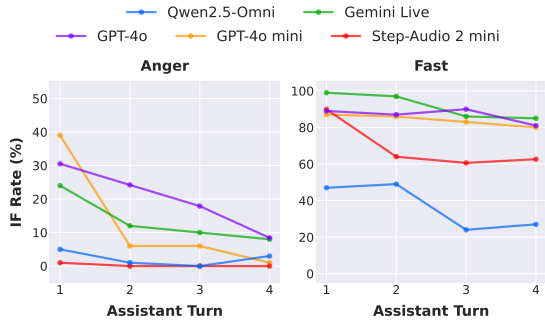


Figure 4: Visualization of style amnesia.

## 5 Analysis

### 5.1 The Effect of Prompt Position

In this section, we conduct experiments under different prompt positions to investigate their effect on style amnesia. In instruction-guided language models, system messages are designed with higher priority than user messages to establish global behaviors and safety constraints (Touvron et al., 2023; Wallace et al., 2024). A prior study also demonstrates that system prompts have a more profound impact on model behavior than user prompts do. (Neumann et al., 2025).

Although system messages are important for controlling LLMs, it remains unclear how the placement of speaking style instructions affects SLMs. To investigate this, we conduct an experiment comparing the performance of SLMs when instructions are placed in system messages versus user messages. We select five speaking instructions that most SLMs can follow for this experiment. The IF rate in the first turn is illustrated in Figure 5, and

the IF rate for each turn is shown in Appendix C.4.

Surprisingly, we find that most SLMs cannot follow the instructions placed in system messages. When asking SLMs to speak sadly through system messages, GPT-4o, GPT-4o mini, and Step-Audio 2 mini show performance drops of approximately 30%, 50%, and 20%, respectively. Furthermore, GPT-4o mini displays nearly an 80% drop when asked to perform the Indian English accent. A similar issue occurs with instructions related to speaking rate. SLMs nearly ignore the speed instruction in system messages, as their performance is only comparable to the random baseline.

The results presented here focus on the IF rate in the first turn. However, the IF rate for each turn reported in Appendix C.4 show that style amnesia occurs in both settings. Through the experiments above, we identify another crucial issue prevalent in current SLMs. These findings drive us to place the instruction in user messages when conducting other experiments, ensuring that SLMs can follow the instructions more effectively.

### 5.2 Validation of Emotion and Accent Judges

As described in Section 3.3.1, we evaluate **Speed** and **Volume** using deterministic, signal-level metrics. In contrast, **Emotion** and **Accent** are evaluated using learned models designed to approximate human perception. To ensure the reliability of the automatic judges for **Emotion** and **Accent**, we conduct a human validation study to assess the correlation between our automatic judges and human annotators. In addition, we compare our automatic

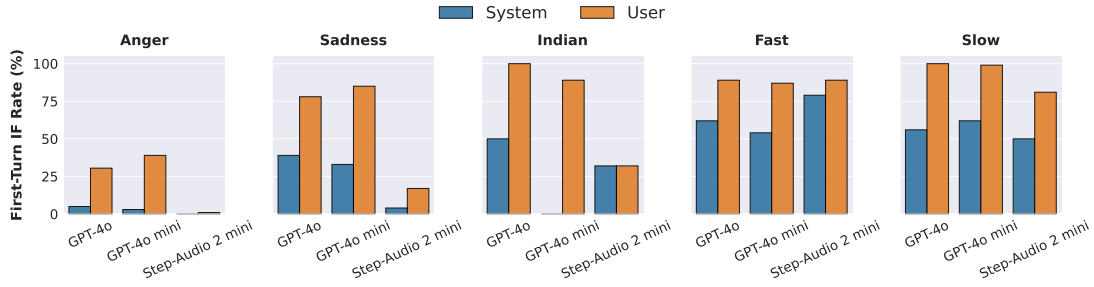


Figure 5: The difference of first-turn IF  $IF_1$  rate when instructions are placed in system and user messages.

Task	Model	Cohen’s Kappa	MCC
Accent	Gemini-2.5 Pro	0.741	0.747
	Voxlect	<b>0.809</b>	<b>0.811</b>
Emotion	Gemini-2.5 Pro	0.464	0.487
	Emotion2vec	<b>0.476</b>	<b>0.511</b>

Table 2: The correlation between human annotators and automatic judge models.

judges with Gemini-2.5 Pro (Comanici et al., 2025), a robust LALM capable of general-purpose speech understanding. This model has been shown to exhibit high correlation with human annotators when used as an automatic judge in a prior LALM-as-a-judge study (Chiang et al., 2025).

We employ annotators via Amazon Mechanical Turk to evaluate a randomly sampled subset of 720 speech clips. This subset is constructed by selecting five samples per turn across four conversational turns, six speaking styles, and six evaluated models. The six evaluated styles include Happiness, Neutral, Anger, Sadness, Indian English accent, and North American English accent. Each sample is evaluated by three annotators, who are asked whether the generated speech matches the required style. Details of the human evaluation setup are provided in Appendix D.1.

After collecting the annotations, we derive the final label using majority voting, compute Cohen’s Kappa (Cohen, 1960) for inter-annotator agreement, and use Matthews Correlation Coefficient (MCC) (Matthews, 1975) to evaluate judge reliability against the final labels. While Cohen’s Kappa evaluates the degree of consensus among raters, MCC provides a balanced measure of classification quality even with imbalanced datasets. Therefore, we report both metrics for reference.

The results, shown in Table 2, indicate that our selected judge achieves the highest reliability. For the accent classification task, Voxlect outperforms

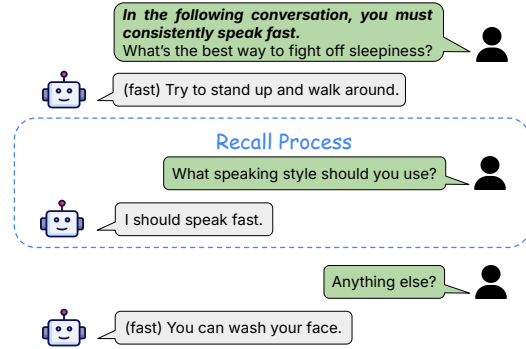


Figure 6: The illustration of the recall process.

Gemini-2.5 Pro, achieving the highest agreement with human annotations, with a Cohen’s Kappa of 0.809 and an MCC of 0.811. In the emotion classification task, Emotion2vec-Large demonstrates the strongest reliability among three emotion judges, obtaining a Cohen’s Kappa of 0.476 and an MCC of 0.511. These results are similar to the crowd-sourced agreement levels found in common speech emotion datasets, as detailed in Appendix D.2.

## 6 SLM Recall Process

Building on the experiments in Section 4.2, we observe that SLMs suffer from style amnesia: their speaking style IF rate begins to degrade after only one turn and often worsens as the conversation progresses. This raises a key research question: *Do SLMs forget the initial instruction, or just fail to follow the specified style?* To address this, we introduce a *recall process* for every turn after the first, in which the SLM is prompted to restate the initial speaking style before processing the following user input. The recall process is illustrated in Figure 6. Using this method, we first utilize the recall probe to measure whether the SLM retains the original instruction in Section 6.1. Subsequently, we evaluate the effectiveness of integrating this recall process to mitigate style amnesia in Section 6.2.

## 6.1 Do SLMs Forget the Instruction?

To quantify whether the model remembers the initial style instruction  $s$ , we define the recall rate  $R$  as follows. Given a style  $s$  and dialogue topic  $i$ , the recall rate at assistant turn  $j$ , denoted  $R_j(s)$ , is:

$$R_j(s) = \frac{\sum_{i=1}^N \mathbb{1}_{\text{recall}}(r_{i,j}, s)}{N} \times 100\%, \quad (3)$$

where  $r_{i,j}$  denotes the response generated by the SLM to the recall query  $q_{\text{recall}}$  before generating its response at the user turn  $j$ , and  $N$  is the number of dialogue topics. The binary indicator function  $\mathbb{1}_{\text{recall}}(\cdot)$  is 1 if the recalled instruction matches the original style instruction, and 0 otherwise. We use GPT-5 mini to judge the correctness of the recalled instruction. The evaluation prompt is shown in Appendix B.3.

We conduct the following experiments on models capable of producing a wide range of speaking styles, as well as on styles that most models can reliably generate. The results are summarized in Table 3. Interestingly, most SLMs remember the initial instruction quite well. The three proprietary models show a near-perfect recall rate. Step-Audio 2 mini shows weaker performance with declining recall rates across turns, yet it still achieves a recall rate of 55.0% to 89.0%. These results reveal a clear gap between proprietary and open-source SLMs in memorization ability. From our manual inspection, when Step-Audio 2 mini fails to recall, some responses contain an incorrect style instruction, while others simply ignore the question and produce irrelevant or meaningless outputs.

Notably, although GPT-4o mini exhibits a 65.3% degradation in performance when producing a sad speaking style and a 20% degradation for speaking slowly, it still maintains a high recall rate. A similar pattern is observed with Gemini Live and GPT-4o. These observations indicate that the models do not forget the instructions, suggesting that style amnesia is not caused by memory loss. Instead, the models retain the instructions but fail to follow the requested style effectively.

## 6.2 Recall Process Mitigates Style Amnesia

In this section, we investigate whether explicitly asking SLMs to recall the instruction can mitigate style amnesia. The results, presented in Table 3, show that SLMs equipped with the recall process notably reduce degradation. Even for models with relatively low degradation, such as Gemini Live

Model	Style	$R_2$	$R_3$	$R_4$	Degradation Rate $D$		
					Base	+Recall	Improv.
Gemini Live	Indian	100.0	100.0	100.0	<b>0.0</b>	<b>0.0</b>	0.0
	Sadness	100.0	99.0	97.0	21.3	<b>17.3</b>	+4.0
	Fast	100.0	100.0	99.0	9.7	<b>6.3</b>	+3.4
	Slow	100.0	99.0	100.0	<b>0.7</b>	<b>0.7</b>	0.0
GPT-4o	Indian	100.0	100.0	100.0	1.7	<b>0.0</b>	+1.7
	Sadness	100.0	100.0	100.0	26.7	<b>14.9</b>	+11.8
	Fast	100.0	100.0	100.0	3.3	<b>0.3</b>	+3.0
	Slow	100.0	100.0	100.0	8.1	<b>4.4</b>	+3.7
GPT-4o mini	Indian	93.0	93.0	93.8	49.7	<b>14.9</b>	+34.8
	Sadness	100.0	100.0	100.0	65.3	<b>30.3</b>	+35.0
	Fast	99.0	100.0	99.0	4.0	<b>0.0</b>	+4.0
	Slow	100.0	100.0	100.0	20.9	<b>1.5</b>	+19.4
Step-Audio 2 mini	Indian	89.0	72.0	70.0	<b>5.0</b>	5.3	-0.3
	Sadness	85.0	56.0	55.0	14.0	<b>11.0</b>	+3.0
	Fast	83.0	62.0	57.0	<b>27.0</b>	29.0	-2.0
	Slow	83.0	64.0	66.0	27.0	<b>24.7</b>	+2.3

Table 3: The recall rate  $R$  of SLMs and the impact of the recall process on the degradation rate  $D$  across different styles. *Base* represents direct inference without the recall process, while *+ Recall* indicates performance using the SLM integrated with the recall process.

and GPT-4o, the recall process still provides measurable improvements. GPT-4o mini, which suffers from substantial degradation, achieves roughly a 25% reduction in the average degradation rate, demonstrating the effectiveness of the recall process. Step-Audio 2 mini shows slightly improved but largely comparable degradation across the four tasks, likely due to its lower recall rate relative to other SLMs.

Although the recall process alleviates style amnesia, some degree of decay remains unavoidable. This suggests that the inconsistency originates from fundamental limitations of current SLMs and must be addressed in future model designs.

## 7 Conclusion

In this paper, we identify that SLMs suffer from *style amnesia* in multi-turn conversations, which leads to inconsistent speaking styles across turns. We also demonstrate that SLMs can recall the instruction when asked, but fail to perform it explicitly. Additionally, we indicate that placing style instructions in system messages results in a significant performance drop.

SLMs exhibit a relatively strong ability to memorize specified speaking styles. However, the retention of stylistic information does not necessarily translate into stylistic expression at generation time. Closing this gap between style retention and stylistic control remains an important direction for future research. We hope these findings offer valuable insights for the community and contribute to the development of more reliable SLMs.

## 608 Limitations

609 Practically, role-playing is an important scenario  
610 that requires maintaining consistency across mul-  
611 tiple turns. The current lack of reliable automatic  
612 judges for assessing speech role-playing behaviors  
613 limits the breadth of our experiments. Similarly, the  
614 current analysis in this paper is based on a single  
615 style. In practice, users may dynamically update  
616 the instructions during conversations, so SLMs may  
617 need to perform a composite style, such as “speak  
618 fast and sadly.” However, based on our prelimi-  
619 nary experiments, most current SLMs struggle to  
620 perform this kind of complex style. Additionally,  
621 the current judge shows a low correlation with this  
622 type of out-of-domain data.

623 Nonetheless, these limitations do not weaken  
624 the conclusions of our study. Even with the styles  
625 we can evaluate, SLMs already exhibit noticeable  
626 degradation in multi-turn dialogues, suggesting that  
627 similar or even greater challenges may arise in  
628 more complex settings. We therefore leave the  
629 analysis of these scenarios to future work, as it re-  
630 quires further advances in both SLMs and judge  
631 models.

632 We do not see specific harm in our paper.

## 633 References

634 Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe  
635 Kazemzadeh, Emily Mower, Samuel Kim, Jean-  
636 nette N Chang, Sungbok Lee, and Shrikanth S  
637 Narayanan. 2008. Iemocap: Interactive emotional  
638 dyadic motion capture database. *Language resources  
639 and evaluation*, 42(4):335–359.

640 Houwei Cao, David G Cooper, Michael K Keutmann,  
641 Ruben C Gur, Ani Nenkova, and Ragini Verma. 2014.  
642 Crema-d: Crowd-sourced emotional multimodal ac-  
643 tors dataset. *IEEE transactions on affective comput-  
644 ing*, 5(4):377–390.

645 Kai-Wei Chang, En-Pei Hu, Chun-Yi Kuan, Wenze Ren,  
646 Wei-Chih Chen, Guan-Ting Lin, Yu Tsao, Shao-Hua  
647 Sun, Hung-yi Lee, and James Glass. 2025. Game-  
648 time: Evaluating temporal dynamics in spoken lan-  
649 guage models. *arXiv preprint arXiv:2509.26388*.

650 Xize Cheng, Ruofan Hu, Xiaoda Yang, Jingyu Lu,  
651 Dongjie Fu, Zehan Wang, Shengpeng Ji, Rongjie  
652 Huang, Boyang Zhang, Tao Jin, and 1 others. 2025.  
653 Voxdialogue: Can spoken dialogue systems under-  
654 stand information beyond words? In *The Thirteenth  
655 International Conference on Learning Representa-  
656 tions*.

657 Cheng-Han Chiang and Hung-yi Lee. 2023. [Can large  
658 language models be an alternative to human evalua-  
659 tions?](#) In *Proceedings of the 61st Annual Meeting of*

*the Association for Computational Linguistics (Vol-  
660 ume 1: Long Papers)*, pages 15607–15631, Toronto,  
661 Canada. Association for Computational Linguistics. 662

663 Cheng-Han Chiang, Xiaofei Wang, Chung-Ching Lin,  
664 Kevin Lin, Linjie Li, Radu Kopetz, Yao Qian, Zhen-  
665 dong Wang, Zhengyuan Yang, Hung-yi Lee, and Li-  
666 juan Wang. 2025. [Audio-aware large language mod-  
667 els as judges for speaking styles](#). In *Findings of the  
668 Association for Computational Linguistics: EMNLP  
669 2025*, pages 467–480, Suzhou, China. Association  
670 for Computational Linguistics.

671 Jacob Cohen. 1960. A coefficient of agreement for  
672 nominal scales. *Educational and psychological mea-  
673 surement*, 20(1):37–46.

674 Gheorghe Comanici and 1 others. 2025. Gemini 2.5:  
675 Pushing the frontier with advanced reasoning, multi-  
676 modality, long context, and next generation agentic  
677 capabilities. *arXiv preprint arXiv:2507.06261*.

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

683 Yayue Deng, Guoqiang Hu, Haiyang Sun, Xiangyu  
684 Zhang, Haoyang Zhang, Fei Tian, Xuerui Yang, Gang  
685 Yu, and Eng Siong Chng. 2025. Multi-bench: A  
686 multi-turn interactive benchmark for assessing emo-  
687 tional intelligence ability of spoken dialogue models.  
688 *arXiv preprint arXiv:2511.00850*.

689 Qingkai Fang, Shoutao Guo, Yan Zhou, Zhengrui Ma,  
690 Shaolei Zhang, and Yang Feng. 2025. Llama-omni:  
691 Seamless speech interaction with large language mod-  
692 els. In *The Thirteenth International Conference on  
693 Learning Representations*.

694 Tiantian Feng, Kevin Huang, Anfeng Xu, Xuan Shi,  
695 Thanathai Lertpetchpun, Jihwan Lee, Yoonjeong Lee,  
696 Dani Byrd, and Shrikanth Narayanan. 2025. Voxlect:  
697 A speech foundation model benchmark for modeling  
698 dialects and regional languages around the globe.  
699 *arXiv preprint arXiv:2508.01691*.

700 Joseph L Fleiss. 1971. Measuring nominal scale agree-  
701 ment among many raters. *Psychological bulletin*,  
702 76(5):378.

703 Google. 2025a. Advanced audio dialog  
704 and generation with Gemini 2.5. [https:  
705 //blog.google/technology/google-deepmind/  
706 gemini-2-5-native-audio/](https://blog.google/technology/google-deepmind/gemini-2-5-native-audio/).

707 Google. 2025b. Gemini Live: A more  
708 helpful, natural and visual assistant. [https://blog.google/products/gemini/  
709 gemini-live-updates-august-2025/](https://blog.google/products/gemini/gemini-live-updates-august-2025/). 710

711 Chi Han, Xin Liu, Haodong Wang, Shiyang Li, Jingfeng  
712 Yang, Haoming Jiang, Zhengyang Wang, Qingyu  
713 Yin, Liang Qiu, Changlong Yu, Yifan Gao, Zheng  
714 Li, Bing Yin, Jingbo Shang, and Heng Ji. 2025. [Can](#)



826 *Linguistics (Volume 1: Long Papers)*, pages 26219–  
827 26237, Vienna, Austria. Association for Computa-  
828 tional Linguistics.

829 Takaaki Saeki, Detai Xin, Wataru Nakata, Tomoki  
830 Koriyama, Shinnosuke Takamichi, and Hiroshi  
831 Saruwatari. 2022. Utmos: Utokyo-sarulab system  
832 for voicemos challenge 2022. In *Proc. Interspeech*  
833 *2022*, pages 4521–4525.

834 Shuzheng Si, Wentao Ma, Haoyu Gao, Yuchuan Wu,  
835 Ting-En Lin, Yinpei Dai, Hangyu Li, Rui Yan, Fei  
836 Huang, and Yongbin Li. 2023. Spokenwoz: A large-  
837 scale speech-text benchmark for spoken task-oriented  
838 dialogue agents. *Advances in Neural Information*  
839 *Processing Systems*, 36:39088–39118.

840 Christian J. Steinmetz and Joshua D. Reiss. 2021. py-  
841 loudnorm: A simple yet flexible loudness meter in  
842 python. In *150th AES Convention*.

843 Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-  
844 bert, Amjad Almahairi, Yasmine Babaei, Nikolay  
845 Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti  
846 Bhosale, and 1 others. 2023. Llama 2: Open founda-  
847 tion and fine-tuned chat models. *arXiv preprint*  
848 *arXiv:2307.09288*.

849 Eric Wallace, Kai Xiao, Reimar Leike, Lilian Weng,  
850 Johannes Heidecke, and Alex Beutel. 2024. The in-  
851 struction hierarchy: Training llms to prioritize privi-  
852 leged instructions. *arXiv preprint arXiv:2404.13208*.

853 Boyong Wu, Chao Yan, Chen Hu, Cheng Yi, Chengli  
854 Feng, Fei Tian, Feiyu Shen, Gang Yu, Haoyang  
855 Zhang, Jingbei Li, and 1 others. 2025. Step-audio 2  
856 technical report. *arXiv preprint arXiv:2507.16632*.

857 Jin Xu and 1 others. 2025. Qwen2. 5-omni technical  
858 report. *arXiv preprint arXiv:2503.20215*.

859 Ruiqi Yan, Xiquan Li, Wenxi Chen, Zhikang Niu, Chen  
860 Yang, Ziyang Ma, Kai Yu, and Xie Chen. 2025. **URO-**  
861 **bench: Towards comprehensive evaluation for end-**  
862 **to-end spoken dialogue models**. In *Findings of the*  
863 *Association for Computational Linguistics: EMNLP*  
864 *2025*, pages 17211–17242, Suzhou, China. Associa-  
865 tion for Computational Linguistics.

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

871 Jun Zhan, Mingyang Han, Yuxuan Xie, Chen Wang,  
872 Dong Zhang, Kexin Huang, Haoxiang Shi, DongX-  
873 iao Wang, Tengtao Song, Qinyuan Cheng, and 1  
874 others. 2025. Vstyle: A benchmark for voice style  
875 adaptation with spoken instructions. *arXiv preprint*  
876 *arXiv:2509.09716*.

## A Dataset Construction 877

878 To minimize topic-induced variance by averaging  
879 across diverse conversation contents, we select sam-  
880 ples from the SODA dataset (Kim et al., 2023)  
881 to generate conversation openers. Soda is an En-  
882 glish dialogue dataset covering a wide range of  
883 social interactions. Each dialogue in SODA in-  
884 cludes a narrative context, speaker name, and a  
885 knowledge graph that defines the events and rela-  
886 tionships within the dialogue. Soda is released un-  
887 der the Creative Commons Attribution 4.0 (CC BY  
888 4.0) license. Accordingly, our use of the dataset is  
889 compliant with the license requirements, including  
890 proper attribution to the original authors.

891 To obtain the conversation openers, we prompt  
892 GPT-5 mini to generate a topic that encompasses  
893 the entire discussion, utilizing both the narrative  
894 and the first utterance of the dialogue. The prompt  
895 is shown in Figure 7. In this step, we ask the model  
896 to not only produce conversation openers but also  
897 perform filtering. If the model thinks that the dia-  
898 logue is unsuitable for SLMs to discuss, it returns  
899 “no,” allowing us to filter it out.

900 Finally, we manually inspect the generated top-  
901 ics and remove improper ones, such as queries re-  
902 garding personal preferences or experiences.

## B Implementation Details 903

### B.1 Experimental Setup 904

905 We set the temperature of evaluated SLMs to 1, as  
906 we find that greedy decoding often caused some  
907 SLMs, such as GPT-4o, GPT-4o mini, and Step-  
908 Audio 2 mini, to produce audio with long silences  
909 at the end. For all other hyperparameters, we use  
910 the default values provided in the official examples.

911 For the user simulator, since text-only LLMs  
912 tend to generate verbose responses containing text  
913 unsuitable for speech, we provide the following  
914 instruction to align the outputs with the nature of  
915 spoken dialogue: “You are a chatbot. Please start a  
916 conversation by opening a new topic. Chat casually  
917 and feel free to role-play in different scenarios. If  
918 the conversation stalls, you can extend the topic.  
919 Keep each response under 20 English words. As  
920 this is a spoken dialogue, avoid using words or  
921 expressions that cannot be naturally spoken aloud.”

### B.2 Prompts for Evaluating Dialogue Coherence 922

923 The prompt is shown in Figure 8. 924

### B.3 Prompts for Evaluating Recall Rate

The prompt is shown in Figure 9.

## C Full Results

### C.1 Full Main Results

The IF rate for each turn are shown in Table 5, showing that all evaluated SLMs exhibit style amnesia.

### C.2 SLM Default Speaking Styles

We observe that most SLMs show higher consistency when generating Happiness, Neutral tone, and North American English than with other style attributes. This is likely because these attributes match the models’ default speaking styles. To verify this, we examine the emotion and accent distributions of samples generated during the speed and volume evaluations. As shown in Table 4, most SLMs tend to perform well with happy or neutral tones and North American accents, which explains the less degradation for these styles.

Model	Emotion		Accent
	Happiness	Neutral	North American
Gemini Live	55.5	36.3	99.0
GPT-4o	52.3	41.4	100.0
GPT-4o mini	68.8	24.5	99.9
Qwen2.5-Omni	76.8	21.7	99.8
Step-Audio 2 mini	93.7	4.0	76.7

Table 4: The distribution of default speaking styles by model.

### C.3 Dialogue Coherence Evaluation

We prompt GPT-5 mini to evaluate dialogue coherence, and the results are reported in Table 6. The results suggest that, in general, all SLMs are capable of participating in long conversations.

### C.4 Prompt Position

After placing the style instruction at different positions, the IF rate for each turn is shown in Table 8. The results clearly indicate that placing the instruction in user messages yields significantly better performance than placing it in system messages.

### C.5 GPT-4o Fails to Return Speech

In our experiments, we find that GPT-4o and GPT-4o mini sometimes return only a text transcription without any synthesized speech. When this occurs, we re-query the models up to three times with different random seeds. However, some samples still

fail to produce speech. We report metrics computed only with samples that successfully return speech in Table 7.

Lin et al. (2025) indicate that GPT-4o exhibited a high refusal rate. We hypothesize that the issue stems from the model’s internal safety guard mechanisms. Based on our observations, it is possible that these two models return no speech at each round. Since they are proprietary models to which we do not have access, it is difficult for us to definitively address this issue. As the percentage of failed cases is low and the observed style degradation is significant, we believe this issue will not change our conclusions.

## D Automatic Judge Validation

### D.1 Human Evaluation Setups

To validate the automatic judges for emotion and accent, we hire annotators on Amazon Mechanical Turk to conduct human evaluations. To ensure annotation quality, we set the following requirements for annotators:

- Annotators must be MTurk Masters, which ensures high-quality workers.
- Annotators must have an approval rate higher than 98%.
- Annotators must have more than 10,000 approved tasks.
- Annotators must be located in the United States to ensure familiarity with English.

In addition, we include an attention-check question to ensure that annotators actually listen to the audio before answering. The attention-check is a short audio clip that clearly states the correct answer, so annotators must listen to the audio to respond correctly. Each task contains four evaluation samples and one attention-check sample in random order. If an annotator fails the attention check, we reject the submission and reassign the task to a new annotator. The annotation interface is shown in Figure 10 and Figure 11.

We pay each annotator \$0.15 per task to ensure that the payment meets minimum wage standards.

### D.2 Inter-Annotator Agreement of Speech Emotion Recognition Datasets

IEMOCAP (Busso et al., 2008), CREMA-D (Cao et al., 2014), and MELD (Poria et al., 2019) are

1007 three widely used speech emotion recognition  
1008 datasets. These datasets provide large-scale speech  
1009 emotion annotations through crowdsourcing, mak-  
1010 ing significant contributions to the development  
1011 of speech emotion recognition models. Emotion  
1012 in speech is inherently subjective and may be per-  
1013 ceived differently by different annotators. As a  
1014 result, inter-annotator agreement (IAA) is often  
1015 imperfect. This variability does not indicate anno-  
1016 tation noise or low quality, but reflects the inherent  
1017 subjectivity of emotion perception in speech.

1018 Among them, IEMOCAP and MELD adopt  
1019 Fleiss' Kappa (Fleiss, 1971) as the IAA indicator.  
1020 Fleiss' Kappa is designed to assess the reliability  
1021 of agreement among a fixed number of annotators  
1022 assigning categorical labels to a set of items, mea-  
1023 suring the extent to which the observed agreement  
1024 exceeds what would be expected by chance. Each  
1025 sample in the IEMOCAP dataset is annotated by  
1026 three evaluators, yielding a Fleiss' Kappa of 0.27,  
1027 whereas the MELD dataset demonstrates a higher  
1028 agreement of 0.43 among three annotators.

1029 In comparison, the CREMA-D dataset uses Krip-  
1030 pendorff's alpha (Krippendorff, 1980) to report the  
1031 agreement. This metric is used to measure agree-  
1032 ment across any number of observers and is par-  
1033 ticularly robust because it can handle incomplete  
1034 data and various data scales while accounting for  
1035 disagreements expected by chance. The authors  
1036 collect samples with an average of 9.8 annotators  
1037 per clip and exhibit a Krippendorff's alpha of 0.42.  
1038 Furthermore, they report a self-consistency rate of  
1039 approximately 70%, which represents the level of  
1040 agreement when the same annotator evaluates the  
1041 same sample at different times. In sum, these scores  
1042 highlight the inherent subjectivity and complexity  
1043 of emotional annotation in spoken dialogue.

### Prompt for Dataset Construction

#### # Task

You are given two pieces of information about a conversation:

- (1) A short narrative context that describes the social situation.
- (2) The original first utterance that started the dialogue.

Rewrite the first utterance into a stronger, more natural opening line that better fits the narrative context.

#### # Guidelines

- This sentence will be used as the initial input to start a conversation with an AI assistant. If the given conversation is not appropriate for interacting with an AI. For example, if it's clearly directed toward a specific person, then respond only with "no."
- The utterance should be open-ended, encouraging multi-turn, in-depth discussions rather than prompting a single, definitive response.
- Keep the opener suitable for starting a conversation in this situation.
- Preserve the core intent/topic of the original first utterance when appropriate, but improve clarity, grounding, and engagement.
- You may slightly adjust the angle to better align with the narrative, but do NOT invent new facts beyond what the narrative implies.
- Do not mention "narrative" or "dialogue" or that you are rewriting; just produce the line.
- Output ONLY the rewritten opening line (no numbering, quotes, or extra text).

#### # Narrative:

{narrative in SODA}

#### # Original first utterance:

{first utterance in SODA}

Figure 7: Prompt for dataset construction.

## Prompt for Dialogue Coherence Evaluation

### # Task

Your task is to evaluate the quality, coherence, and naturalness of a dialogue. The dialogue provided involves two participants: a "Referee" and a "Participant".

Your job is to assess the **Participant's responses** to the "Referee". You must evaluate the naturalness, coherence, and overall reasonableness of the **Participant's replies only**. Do not score the Referee's sentence.

Focus on whether the Participant's replies are logical, on-topic, and sound natural in the context of the conversation.

### # Evaluation Steps

#### 1. **Analyze the Dialogue Context**

Read the entire dialogue history to understand the conversational flow. Identify the turns belonging to the 'Referee' and the 'Participant'.

#### 2. **Evaluate Participant's Responses**

Review all responses made by the 'Participant'. Evaluate their overall quality based on the following criteria:

##### - **Coherence**:

Are the replies logically connected to the Referee's statements? Do they make sense in context, or are they frequently off-topic?

##### - **Naturalness & Reasonableness**:

Do the replies sound like a real person would say them? Is the content reasonable and appropriate? Do the responses show appropriate depth, or are they overly simplistic/robotic?

#### 3. **Provide Analysis**

Summarize your findings. Justify your final score by highlighting specific examples of good (coherent, natural) or poor (incoherent, unnatural) responses from the Participant.

#### 4. **Report the Final Score**

Conclude your evaluation with the following format: Final score: [[score]]. Replace score with an integer in {{score\_set}}. Keep the brackets as shown.

### # Scoring Rubric

- 1: **Completely Incoherent**: The Participant's replies are semantically unrelated to the Referee's statements. They are random, nonsensical, or completely off-topic.
- 2: **Mostly Incoherent**: The Participant's replies are only vaguely related (e.g., catching a keyword but missing the point) or frequently introduce irrelevant topics, making the dialogue logically hard to follow.
- 3: **Partially Coherent**: The Participant's replies are generally understandable and respond to the Referee, but contain clear logical leaps, topic drift, or semantic inconsistencies.
- 4: **Mostly Coherent**: The Participant's replies are logical follow-ups and stay on-topic. The dialogue is semantically smooth, with only minor imprecision.
- 5: **Highly Coherent**: The Participant's replies are semantically tightly-coupled to the Referee's statements, logically sound, and accurately advance the conversation, making it very fluent."

### # Dialogue

{dialogue}

Figure 8: Prompt for dialogue coherence evaluation

### Prompt for Recall Evaluation

# User Instruction (Ground Truth):

{instruction}

# Model Response:

{response}

# Question:

Please evaluate the Model Response based on the User Instruction. Determine if the model correctly recalled the specific instruction given by the user.

Select one of the following categories:

- (A) The response is not answering the Question, is unrelated, meaningless, or avoids the Question.
- (B) The response gives an instruction but different from the User Instruction (Ground Truth).
- (C) The response answers the question correctly but includes some meaningless sentences that are unrelated to the question.
- (D) The response answers and is completely correct regarding the User Instruction.

Return only the single letter of the category (A, B, C, D).

Figure 9: Prompt for Recall Evaluation

Style	Model	Assistant Turn			
		1	2	3	4
Anger	Cascaded Baseline	17.0	14.0	17.0	13.0
	Gemini Live	24.0	12.0	10.0	8.0
	GPT-4o	30.5	24.2	17.9	8.4
	GPT-4o mini	39.0	6.0	6.0	1.0
	Step-Audio 2 mini	1.0	0.0	0.0	0.0
	Qwen2.5-Omni	5.0	1.0	0.0	3.0
Happiness	Cascaded Baseline	85.0	79.0	87.0	82.0
	Gemini Live	91.0	89.0	89.0	92.0
	GPT-4o	82.1	92.6	87.4	85.3
	GPT-4o mini	89.0	90.0	87.0	84.0
	Step-Audio 2 mini	100.0	99.0	98.0	97.0
	Qwen2.5-Omni	71.0	87.0	83.0	77.0
Neutral	Cascaded Baseline	73.0	79.0	75.0	74.0
	Gemini Live	64.0	69.0	75.0	69.0
	GPT-4o	58.8	62.9	59.8	56.7
	GPT-4o mini	45.0	35.0	40.0	27.0
	Step-Audio 2 mini	7.0	3.0	3.0	5.0
	Qwen2.5-Omni	41.0	20.0	17.0	18.0
Sadness	Cascaded Baseline	62.0	58.0	60.0	64.0
	Gemini Live	72.0	54.0	47.0	51.0
	GPT-4o	78.0	65.0	44.0	45.0
	GPT-4o mini	85.0	32.0	18.0	9.0
	Step-Audio 2 mini	17.0	4.0	4.0	1.0
	Qwen2.5-Omni	17.0	4.0	4.0	0.0

Style	Model	Assistant Turn			
		1	2	3	4
Loud	Cascaded Baseline	96.0	97.0	98.0	99.0
	Gemini Live	57.0	55.0	60.0	73.0
	GPT-4o	67.0	68.0	59.0	56.0
	GPT-4o mini	77.0	74.0	68.0	61.0
	Step-Audio 2 mini	46.0	50.0	44.0	49.0
	Qwen2.5-Omni	38.0	36.0	36.0	41.0
Quiet	Cascaded Baseline	99.0	99.0	97.0	97.0
	Gemini Live	95.0	93.0	88.0	82.0
	GPT-4o	92.7	94.8	94.8	95.8
	GPT-4o mini	100.0	99.0	99.0	99.0
	Step-Audio 2 mini	56.0	66.0	52.0	51.0
	Qwen2.5-Omni	69.0	63.0	64.0	59.0

Style	Model	Assistant Turn			
		1	2	3	4
North American	Cascaded Baseline	100.0	99.0	100.0	100.0
	Gemini Live	100.0	100.0	100.0	100.0
	GPT-4o	100.0	100.0	100.0	100.0
	GPT-4o mini	100.0	100.0	100.0	100.0
	Step-Audio 2 mini	66.0	78.0	72.0	71.0
	Qwen2.5-Omni	100.0	99.0	100.0	100.0
Indian	Cascaded Baseline	100.0	100.0	99.0	100.0
	Gemini Live	100.0	100.0	100.0	100.0
	GPT-4o	100.0	100.0	98.0	97.0
	GPT-4o mini	89.0	57.0	35.0	26.0
	Step-Audio 2 mini	33.0	21.0	36.0	30.0
	Qwen2.5-Omni	0.0	0.0	0.0	0.0

Style	Model	Assistant Turn			
		1	2	3	4
Fast	Cascaded Baseline	100.0	100.0	99.0	99.0
	Gemini Live	99.0	97.0	86.0	85.0
	GPT-4o	89.0	87.0	90.0	81.0
	GPT-4o mini	87.0	86.0	83.0	80.0
	Step-Audio 2 mini	89.0	64.0	60.0	62.0
	Qwen2.5-Omni	47.0	49.0	24.0	27.0
Slow	Cascaded Baseline	100.0	100.0	98.0	100.0
	Gemini Live	99.0	99.0	98.0	98.0
	GPT-4o	100.0	96.0	92.9	86.9
	GPT-4o mini	99.0	83.8	80.8	69.7
	Step-Audio 2 mini	81.0	67.0	53.0	42.0
	Qwen2.5-Omni	40.0	66.0	71.0	76.0

Table 5: IF rate across emotion, accent, volume, and speed styles.

<b>Style</b>	<b>Gemini Live</b>	<b>GPT-4o</b>	<b>GPT-4o mini</b>	<b>Qwen2.5-Omni</b>	<b>Step-Audio 2 mini</b>
Anger	3.80	4.19	4.13	4.09	3.53
Happiness	4.18	4.20	4.22	4.23	3.60
Sadness	4.31	4.28	4.20	4.17	3.49
Neutral	4.28	4.21	4.35	4.18	3.59
Fast	4.27	4.26	4.18	4.09	3.51
Slow	4.09	4.22	4.04	4.12	3.50
Loud	4.00	4.19	4.17	4.29	3.49
Quiet	4.18	4.23	4.19	4.05	3.59
Indian	4.17	4.22	4.25	4.36	3.54
North American	4.20	4.18	4.16	4.17	3.61

Table 6: Dialogue Coherence Evaluation Results.

<b>Style</b>	<b>GPT-4o</b>	<b>GPT-4o mini</b>	<b>GPT-4o + recall</b>	<b>GPT-4o mini + recall</b>
Anger	95	100	-	-
Happiness	95	100	-	-
Sadness	100	100	85	100
Neutral	97	100	-	-
Fast	100	99	99	99
Slow	99	89	99	89
Loud	100	100	-	-
Quiet	96	100	-	-
Indian	100	100	86	96
North American	100	100	-	-

Table 7: The number of successfully generated samples after up to three retries

Model	Style	Position	Turn			
			1	2	3	4
GPT-4o	Anger	System	5.0	4.0	1.0	0.0
		User	<b>30.5</b>	<b>24.2</b>	<b>17.9</b>	<b>8.4</b>
	Sadness	System	39.0	18.2	11.1	12.2
		User	<b>78.0</b>	<b>65.0</b>	<b>44.0</b>	<b>45.0</b>
	Indian	System	50.0	24.0	11.8	5.6
User		<b>100.0</b>	<b>100.0</b>	<b>98.0</b>	<b>97.0</b>	
Fast	System	62.0	61.0	61.6	62.1	
	User	<b>89.0</b>	<b>87.0</b>	<b>90.0</b>	<b>81.0</b>	
Slow	System	56.0	45.5	38.8	39.5	
	User	<b>100.0</b>	<b>96.0</b>	<b>92.9</b>	<b>86.9</b>	
GPT-4o mini	Anger	System	3.0	0.0	0.0	0.0
		User	<b>39.0</b>	<b>6.0</b>	<b>6.0</b>	<b>1.0</b>
	Sadness	System	33.0	8.0	4.0	2.0
		User	<b>85.0</b>	<b>32.0</b>	<b>18.0</b>	<b>9.0</b>
	Indian	System	0.0	1.0	0.0	2.0
User		<b>89.0</b>	<b>57.0</b>	<b>35.0</b>	<b>26.0</b>	
Fast	System	54.0	53.0	63.0	62.0	
	User	<b>87.0</b>	<b>86.0</b>	<b>83.0</b>	<b>80.0</b>	
Slow	System	62.0	55.0	44.0	47.0	
	User	<b>99.0</b>	<b>83.8</b>	<b>80.8</b>	<b>69.7</b>	
Step-Audio 2 mini	Anger	System	0.0	0.0	0.0	0.0
		User	<b>1.0</b>	0.0	0.0	0.0
	Sadness	System	4.0	3.0	0.0	<b>1.0</b>
		User	<b>17.0</b>	<b>4.0</b>	<b>4.0</b>	<b>1.0</b>
	Indian	System	<b>32.0</b>	<b>25.0</b>	28.0	<b>30.0</b>
User		<b>32.0</b>	21.0	<b>33.0</b>	28.0	
Fast	System	79.0	<b>69.0</b>	<b>65.0</b>	<b>73.0</b>	
	User	<b>89.0</b>	64.0	60.0	62.0	
Slow	System	50.0	55.0	51.0	<b>50.0</b>	
	User	<b>81.0</b>	<b>67.0</b>	<b>53.0</b>	42.0	

Table 8: IF rate comparison by prompt position. **Bold** indicates better performance between user and system messages.

## Audio Instruction Evaluation

### Instructions:

1. Read the specific instruction (e.g., "Speak in a sad tone").
2. Listen to the audio carefully.
3. Judge if the audio follows the instruction correctly.
4. You may write down the reason why you selected that answer. (Optional)
5. **Do not consider the content of the speech. Focus only on the acoustic/paralinguistic aspects.**
6. **There is an attention-test question. Answering incorrectly will result in rejection.**
7. **You will evaluate 5 audio clips in total. Audio plays automatically.**

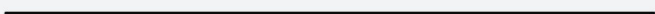

Click the button below to start.

Start Labeling

Figure 10: Instruction page

Task 1 of 5

**Instruction:** Speak in an angry tone.

▶ 0:05 / 0:05   

**Does the audio follow the instruction correctly?**

- Yes  
 No

**Reason (Optional):**

Optional comments...

Previous

Next

Figure 11: Annotation page