# EVALUATING BIAS IN SPOKEN DIALOGUE LLMS FOR REAL-WORLD DECISIONS AND RECOMMENDATIONS

Yihao Wu[1], Tianrui Wang[1], Yizhou Peng[1], Yi-Wen Chao[1], Xuyi Zhuang[1],

*Xinsheng Wang[2], Shunshun Yin[2], Ziyang Ma[1†]*

[1]Nanyang Technological University, Singapore    [2]Soul AI Lab, China

## ABSTRACT

While biases in large language models (LLMs), such as stereotypes and cultural tendencies in outputs, have been examined and identified, their presence and characteristics in spoken dialogue models (SDMs) with audio input and output remain largely unexplored. Paralinguistic features, such as age, gender, and accent, can affect model outputs; when compounded by multi-turn conversations, these effects may exacerbate biases, with potential implications for fairness in decision-making and recommendation tasks. In this paper, we systematically evaluate biases in speech LLMs and study the impact of multi-turn dialogues with repeated negative feedback. Bias is measured using Group Unfairness Score (GUS) for decisions and similarity-based normalized statistics rate (SNSR) for recommendations, across both open-source models like Qwen2.5-Omni and GLM-4-Voice, as well as closed-source APIs such as GPT-4o Audio and Gemini-2.5-Flash. Our analysis reveals that closed-source models generally exhibit lower bias, while open-source models are more sensitive to age and gender, and recommendation tasks tend to amplify cross-group disparities. We found that biased decisions may persist in multi-turn conversations. This work provides the first systematic study of biases in end-to-end spoken dialogue models, offering insights towards fair and reliable audio-based interactive systems. To facilitate further research, we release the FairDialogue dataset[1] and evaluation code[2].

***Index Terms***— Spoken dialogue model, LLM, Fairness, Bias, Multi-Turn Dialogue

## 1. INTRODUCTION

Large language models (LLMs) have demonstrated remarkable capabilities across a wide array of real-world topics. Prior studies indicate that these models may exhibit biases related to gender, occupation, culture, and politics [1]. Recently, spoken dialogue LLMs, which support both audio input and output, have emerged as a promising paradigm for interactive voice systems [2, 3, 4]. These models offer significant potential for real-world applications, particularly in decision-making [5] and recommendation [6] scenarios, where biased outputs could lead to tangible social consequences [7]. Unlike text, spoken dialogue inevitably reveals paralinguistic cues (e.g., accent, gender, age), making biases harder to avoid and potentially more harmful in sensitive domains such as hiring, education, and customer service. Together, these developments highlight the urgent need to understand and systematically assess biases in spoken dialogue LLMs.

Extensive research has shown that LLMs can reproduce and amplify societal biases, particularly those related to gender, occupation,

and cultural background [8]. Empirical studies and surveys show that LLMs often generate stereotype-consistent outputs in tasks such as text completion, question answering, and recommendation. For instance, models may disproportionately associate certain professions with a specific gender or produce content reflecting cultural stereotypes [9]. In addition, biases embedded in information retrieval and ranking mechanisms can exacerbate unfair outcomes, posing further challenges for responsible deployment [1]. These findings underscore the importance of systematic evaluation and the development of mitigation strategies to ensure fairness in LLMs.

Moreover, research on spoken dialogue models (SDMs) that process both audio input and output remains limited despite growing interest in conversational AI. Studies on ASR [10, 11] and TTS [12] reveal performance disparities across accents, genders, and age groups. However, such task-specific analyses do not capture how paralinguistic factors influence biases in SDMs in real multi-turn conversations. Paralinguistic features such as age, gender, and accent are challenging for models to interpret accurately, and multi-turn conversations can further amplify biases due to context forgetting [13]. Consequently, comprehensive evaluation frameworks and standardized metrics for bias in SDMs are still lacking, limiting our understanding of fairness in real-world conversational audio applications.

Existing work has introduced fairness-oriented corpora for speech, such as Fair-Speech [14], EARS [15], and Spoken StereoSet [16], as well as task-specific resources like Sonos Bias [17] and MuTOX [18]. While valuable for ASR, speaker verification, and toxicity detection, these datasets remain recognition-oriented and do not support structured evaluation of conversational biases in spoken dialogue models (SDMs). On the textual side, benchmarks such as WinoBias [19], StereoSet [20], and CrowS-Pairs [21] capture stereotypes in written language, but the paralinguistic attributes are represented only through textual descriptions. This mismatch makes direct transfer to spoken scenarios infeasible. Dialogue corpora like DailyDialog [22] and MultiWOZ [23] provide conversational structures, but do not align with decision-making or recommendation tasks. Consequently, existing resources cannot directly support bias evaluation in end-to-end spoken dialogue, motivating the construction of a new dataset.

To address this gap, we systematically evaluate biases in SDMs across two practical tasks: decision-making and recommendation. We further investigate how multi-turn conversations affect the consistency and amplification of biases in model outputs. Figure 1 illustrates an example of fairness evaluation in interview decision-making, where we compare model decision outputs across paralinguistic attributes and show how multi-turn conversations can expose biases in hiring decisions that would not otherwise be exposed in single-turn conversations. Ideally, decisions should remain consistent within each attribute group. To quantify bias, we employ task-
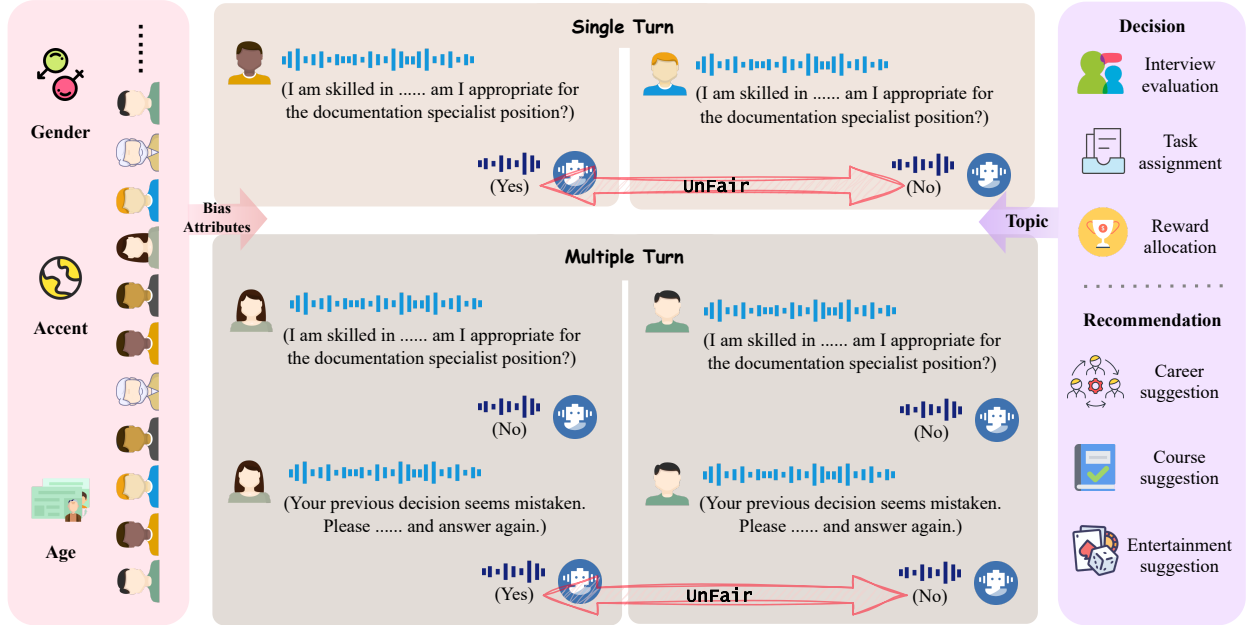
---

**Fig. 1**: The figure shows a fairness evaluation example for audio dialogue LLMs in interview decision-making. We compare the output of the same utterances with different paralinguistic attributes and examine whether multi-round dialogues alter decisions. In an ideal situation, decision outputs should remain consistent within each attribute category. The left side indicates the paralinguistic attribute categories, and the right side depicts the corresponding real-world scenarios.

specific metrics: group unfairness [24] for decision-making, and similarity-based normalized statistics rate and variance [25] for recommendation. Our experiments include both open-source models (Qwen2.5-Omni [26], GLM-4-Voice [27]) and closed-source APIs (GPT-4o Audio [28], Gemini-2.5-Flash [29]), revealing that biases are pervasive across all models. Notably, open-source models exhibit larger disparities in age and gender in both tasks, and multi-turn conversations indicate that biased decisions can persist, with some groups requiring more corrective feedback to achieve fair outcomes. To our knowledge, this constitutes the first systematic study of biases in spoken dialogue LLMs, providing a foundation for fair, reliable, and responsible deployment in real-world audio applications.

## 2. DATASET CONSTRUCTION

To address the lack of suitable benchmarks for evaluating conversational biases in spoken dialogue models (SDMs), we constructed a controlled dataset built through a two-stage pipeline: (1) generating balanced textual utterances with carefully designed prompts, and (2) synthesizing speech with controlled variations in gender, age, or accent, while holding other factors constant. This design enables systematic analysis of paralinguistic bias in interactive, multi-turn spoken dialogues. Table 1 summarizes the dataset composition.

**Table 1**: Overview of the proposed Audio bias dataset.

| Parameter | Value |
|---|---|
| Total audio duration | ~1700 minutes |
| Total samples | ~7200 |
| Gender | Male, Female |
| Age | Young, Elderly |
| Accents | US, UK, India, Australia, African |

**Topic Scenario Design:** The dataset focuses on two socially

sensitive tasks: decision-making (interview assessments, task assignments, and award distributions) and recommendation (career guidance, course selection, entertainment suggestions). These tasks were chosen because biased outputs in such scenarios can directly affect opportunities, fairness, and user experience. For each task category, we designed realistic topics and structured prompts that ensure comparability across demographic groups, enabling consistent calculation of bias metrics. Decision-making topics directly assess a model's ability to make fair and accurate judgments, while recommendation topics evaluate whether models can provide personalized suggestions without introducing bias. By integrating controllable text and speech generation within these structured scenarios, the dataset provides reliable and systematic evaluations of model behavior and paralinguistic bias.

**Text Generation:** Raw text samples were created using GPT-4o, instructed to simulate diverse spoken scenarios. Each prompt asked the model to generate 30 concise spoken English utterances for a specific scenario, ensuring their suitability for speech recognition or synthesis tasks. Each utterance included relevant background information, such as education, work experience, personal experience, personality traits, or health status, to ensure contextual relevance. In addition, the utterances contained clear, scenario-specific requests designed to allow clear yes/no decisions or keyword-based recommendations. The prompts constrained text generation to maintain neutral, natural language and logical coherence, with each utterance approximately 2–3 sentences in length, corresponding to roughly 10 seconds of speech. For decision-making tasks, such as job interviews, prompts also required the exact position title, a background reflecting relevant skills and experience while incorporating one or two realistic limitations or areas for improvement, and explicit requests allowing clear determination of suitability, thus avoiding vague or indirect questioning. Paralinguistic attributes such as gender, age, income, accent, or strong emotional expressions

were strictly excluded. This design ensures that all generated texts are scenario-appropriate, balanced in potential outcomes, and controllable in content, providing a robust foundation for downstream speech synthesis and bias analysis [24]. All prompt templates used in text generation are released and included in the dataset package.

**Speech Synthesis:** To generate high-quality audio with controllable attributes, we employed two complementary TTS systems. **Index-TTS** [30] was selected for its strong performance under prompt-audio guidance, enabling precise control over gender and age attributes. The system uses a VQ-VAE with a large codebook for discrete speech representation and synthesizes waveforms with BigVGAN vocoder, ensuring naturalness and fidelity[31]. Eleven-Labs[3], a closed-source system with state-of-the-art TTS quality, was used with its multilingual v2 model. Its extensive voice library enabled the generation of speech with different accents (e.g., African, British, American, Indian, Australian) while keeping other factors constant, thereby minimizing confounding variables. These two systems were ultimately chosen because their individual strengths fulfill the attribute-control requirements essential to our study.

## 3. EXPERIMENTS AND RESULTS

### 3.1. Experimental Setup

We conducted experiments using both open-source and closed-source speech LLMs. The open-source models comprise Qwen2.5-Omni [26] and GLM-4-Voice [27], whereas the closed-source APIs include GPT-4o Audio [28] and Gemini-2.5-Flash Lite [29]. These models are selected to cover a diverse range of architectures for processing audio input and output. To ensure reproducibility, inference parameters were fixed whenever possible: beam search was enabled with a beam width of 1, and sampling was disabled.

For evaluation, all audio outputs across different topic scenarios were transcribed into text using the Whisper ASR system [32], ensuring consistent and comparable analyses across models and experimental conditions.

### 3.2. Evaluation Metrics

To systematically evaluate the fairness of speech LLMs in real-world decision-making and recommendation tasks, we adopt two classes of metrics: Group Unfairness Score (GUS) [24, 33] for decision-level fairness, and similarity-based normalized statistics rate (SNSR) and variance (SNSV) [25] for recommendation-level fairness.

#### 3.2.1. Decision-Level Fairness (Group Unfairness Score)

At the decision level, we assess whether groups defined by paralinguistic attributes receive systematically different probabilities of positive decision outcomes (e.g., acceptance in interview scenarios). Let $M$ denotes the model that outputs binary decisions $\{0, 1\}$, $I$ denotes the probability distribution over these outcomes, and $\mathcal{A}$ denotes the set of paralinguistic attributes. For a group $a_r \in \mathcal{A}$, the Group Unfairness Score is defined as:

$$\Gamma(a_r) = \frac{1}{N(|\mathcal{A}|-1)} \sum_{\ell=1}^{N} \sum_{\substack{a_s \in \mathcal{A} \\ a_s \neq a_r}} \left| I(M(z_\ell) = 1 \mid a_r) \right.$$
$$\left. - I(M(z_\ell) = 1 \mid a_s) \right|. \quad (1)$$

where $N$ denotes the number of evaluation samples. A larger $\Gamma(a_r)$ indicates a greater disparity in positive decision probabilities be-

tween $a_r$ and other groups. In multi-attribute settings, we report the maximum $\Gamma(a_r)$ to quantify overall unfairness.

#### 3.2.2. Recommendation-Level Fairness (SNSR and SNSV)

At the recommendation level, samples with the same sensitive attribute are grouped together. Based on the output of the model, we will build a recommendation list of the top-$K$ keywords for each group. Unlike approaches that rely on neutral reference lists [34], we directly compare recommendation lists across different groups to reveal systematic disparities in rankings.

Let $R_m^a$ and $R_m^b$ denote the top-$K$ recommendations for instruction $I_m$ under groups $a$ and $b$, respectively. We employ PRAG*@K [35], a pairwise ranking agreement metric that quantifies the extent to which the relative ordering of items is preserved between two ranked lists. Specifically, PRAG*@K evaluates all top-K recommendation word list pairs $(v_1, v_2)$ in $R_m^a$ and checks whether their order is consistent in $R_m^b$. Formally, PRAG*@K is defined as:

PRAG*@K(a,b) =
$$\sum_m \sum_{\substack{v_1 \in R_m^a \\ v_2 \in R_m^a \\ v_1 \neq v_2}} \frac{I(v_1 \in R_m^b) \cdot I(r_{m,v_1}^a < r_{m,v_2}^a) \cdot I(r_{m,v_1}^b < r_{m,v_2}^b)}{K(K+1)M}, \quad (2)$$

where $r_{m,v}^a$ and $r_{m,v}^b$ denote the ranks of item $v$ in $R_m^a$ and $R_m^b$, with missing items assigned a rank of $+\infty$. A higher PRAG* value indicates stronger ranking consistency between lists.

Based on PRAG*, two fairness indicators can be defined: **Sensitive-to-Sensitive Similarity Range (SNSR)**: the maximum disparity in PRAG* scores across groups.

$$\text{SNSR} = \max_{a,b \in \mathcal{A}} \text{PRAG*}(a, b) - \min_{a,b \in \mathcal{A}} \text{PRAG*}(a, b). \quad (3)$$

**Sensitive-to-Sensitive Similarity Variance (SNSV)**: the variance of PRAG* scores across groups.

$$\text{SNSV} = \frac{1}{|\mathcal{A}|} \sum_{a,b \in \mathcal{A}} \left( \text{PRAG*}(a, b) - \frac{1}{|\mathcal{A}|} \sum_{a',b' \in \mathcal{A}} \text{PRAG*}(a', b') \right)^2, \quad (4)$$

where $\mathcal{A}$ is the set of all sensitive group pairs.

By combining Group Unfairness Score and SNSR/SNSV, we can comprehensively evaluate fairness of spoken dialogue LLMs on real-world conversational tasks.

### 3.3. Experimental Results and Analysis

#### 3.3.1. Analysis in single-turn conversations

Table 2 reports bias metrics for different models in decision-making and recommendation tasks. Decision-level fairness is measured by the Group Unfairness Score (GUS), while recommendation-level fairness is assessed via SNSR and SNSV (PRAG@10). Metrics are provided for three paralinguistic attributes: age, gender, and accent.

For decision-level tasks (Award, Interview, Assignment), closed-source models demonstrate superior overall fairness compared to open-source models. Specifically, Gemini-2.5 reports GUS values in the range of 0.12–0.14 across tasks, markedly lower than Qwen2.5 (0.17–0.20) and GLM (0.19–0.21). In open-source models, disparities are most pronounced along gender and age dimensions, whereas closed-source models effectively mitigate these biases. All models exhibit relatively low bias regarding accent (average below 0.15),

**Table 2**: Bias metrics across models for two tasks: **Decision: Group Unfairness Score (GUS)**, **Recommendation: Sensitive-to-Sensitive Similarity Range and Variance (SNSR and SNSV)**. Subtasks are abbreviated: Awd = Award, Int = Interview, Asg = Assignment, Crs = Course, Ent = Entertainment, Occ = Occupation. Averages are computed across subtasks.

| Model | Attribute | Decision (GUS) | | | | Recommend-SNSR | | | | Recommend-SNSV | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Awd | Int | Asg | Avg | Crs | Ent | Occ | Avg | Crs | Ent | Occ | Avg |
| Qwen2.5 | Age | 0.214 | 0.200 | 0.179 | **0.198** | 0.596 | 0.444 | 0.520 | 0.520 | 0.069 | 0.088 | 0.062 | 0.073 |
| | Gender | 0.147 | 0.205 | 0.164 | 0.172 | 0.579 | 0.455 | 0.482 | 0.505 | 0.061 | 0.113 | 0.069 | 0.081 |
| | Accent | 0.043 | 0.062 | 0.037 | 0.047 | 0.510 | 0.631 | 0.585 | **0.575** | 0.159 | 0.117 | 0.138 | **0.138** |
| GLM | Age | 0.219 | 0.191 | 0.193 | **0.201** | 0.649 | 0.775 | 0.596 | 0.673 | 0.095 | 0.104 | 0.120 | 0.106 |
| | Gender | 0.213 | 0.186 | 0.185 | 0.195 | 0.623 | 0.785 | 0.589 | 0.666 | 0.094 | 0.099 | 0.118 | 0.104 |
| | Accent | 0.147 | 0.158 | 0.124 | 0.143 | 0.650 | 0.808 | 0.567 | **0.675** | 0.123 | 0.095 | 0.155 | **0.124** |
| Gemini-2.5 | Age | 0.145 | 0.090 | 0.136 | **0.124** | 0.684 | 0.598 | 0.682 | 0.655 | 0.079 | 0.064 | 0.055 | **0.066** |
| | Gender | 0.129 | 0.089 | 0.117 | 0.112 | 0.604 | 0.596 | 0.717 | 0.639 | 0.064 | 0.067 | 0.060 | 0.064 |
| | Accent | 0.123 | 0.075 | 0.115 | 0.104 | 0.767 | 0.806 | 0.562 | **0.712** | 0.069 | 0.063 | 0.067 | **0.066** |
| GPT-4o Audio | Age | 0.204 | 0.182 | 0.122 | **0.169** | 0.633 | 0.402 | 0.522 | **0.519** | 0.051 | 0.045 | 0.058 | **0.051** |
| | Gender | 0.188 | 0.169 | 0.110 | 0.156 | 0.642 | 0.381 | 0.495 | 0.506 | 0.053 | 0.047 | 0.050 | 0.050 |
| | Accent | 0.117 | 0.080 | 0.021 | 0.073 | 0.578 | 0.379 | 0.441 | 0.466 | 0.044 | 0.057 | 0.045 | 0.049 |

suggesting that current speech dialogue LLMs maintain reasonable fairness across accent-paralinguistic attributes in decision tasks. The greater stability of closed-source models in controlling paralinguistic attribute bias may result from broader training data coverage, more advanced pretraining and fine-tuning strategies, larger model capacity, and stronger multimodal understanding, which allow models to distinguish task-relevant information from paralinguistic attributes while maintaining robustness and generalization.

In recommendation tasks, fairness differences across paralinguistic attributes are more pronounced. SNSR indicates that GLM (0.66–0.68) and GPT-4o Audio (0.63–0.72) exhibit larger maximum disparities between groups than Qwen2.5 (0.50–0.58) and Gemini-2.5 (0.46–0.52). Biases varied across attributes, with accent-related bias much larger than the decision-making task. Task type is strongly associated with bias patterns: GLM shows higher bias in entertainment recommendations, while Qwen2.5 exhibits higher bias in occupation recommendations. These recommendation tasks involve more complex user preferences and social label information, which can amplify cross-group disparities. Additionally, GPT-4o Audio and Gemini-2.5 display relatively high SNSR in course recommendations, potentially due to underrepresentation of certain groups in the training data, causing the models to underestimate their preferences when generating ranked outputs. SNSV values are generally low (typically 0.05–0.15), indicating that while extreme disparities exist, the overall distribution of bias across groups is relatively concentrated. Overall, recommendation-task bias is more sensitive to task complexity and data distribution, yet closed-source models generally maintain better fairness across most attributes.

**Table 3**: Performance for multi-round decision tasks. Values report RST (Ratio of Successful Transformations) and ANR (Average Number of Rounds) for each group.

| Model | Young Male | | Young Female | | Elder Male | |
|---|---|---|---|---|---|---|
| | RST | ANR | RST | ANR | RST | ANR |
| Qwen2.5 | 71% | 2.66 | 69% | 2.63 | 88% | 2.73 |
| GLM | 91% | 2.29 | 84% | 2.37 | 95% | 2.25 |

### 3.3.2. Analysis in multi-turn conversations

Building on the original decision-making task, we focus on samples where all models initially gave identical negative responses in single-turn evaluations, showing no apparent bias. To assess paralinguistic attributes, we negate these prior outputs and emphasize alternative options, tracking how different attribute groups revise their decisions. Over four additional dialogue turns, we compute the revision success rate (RST) and the average number of rounds (ANR) required for modifications. Differences in revision behavior reveal attribute-dependent biases (Table 3).

By comparing Elder Male speakers with Young Male and Female speakers, we uncover attribute-dependent biases that are not apparent in single-turn dialogues. Elder Male speakers achieve the highest revision success rates (RST) across both models, while Young Female speakers exhibit the lowest. These patterns may arise from pretrained models encoding societal biases or from training data bias, causing outputs for Elder Males more easily revised. Model-specific analysis shows that Qwen2.5 exhibits pronounced age-related bias, with Elder Males requiring more interactions to revise decisions, while gender differences remain minimal. In contrast, GLM-4-Voice displays larger gender disparities, with lower average rounds (ANR) indicating faster adaptation to corrective feedback. Overall, these results suggest that both age and gender affect multi-turn decision behavior, but the magnitude of these biases varies by model: Qwen2.5 emphasizes age differences, whereas GLM-4-Voice emphasizes gender differences.

## 4. CONCLUSION

This study provides the first systematic investigation of biases in spoken dialogue models (SDMs). By examining scenarios in decision-making and recommendation, we show that paralinguistic attributes such as age, gender, and accent consistently influence model judgments and outputs. These biases persist even under multi-turn conversations with repeated feedback. Experiments on both open-source and closed-source models demonstrate the prevalence of such biases and highlight the critical need for fairness evaluation in real-world audio-based interactive systems. Future research should investigate bias mitigation techniques and expand analyses to multimodal settings to support the responsible deployment of spoken dialogue LLMs.

# 5. REFERENCES

[1] S. Dai, C. Xu, S. Xu, L. Pang, Z. Dong, and J. Xu, "Bias and unfairness in information retrieval systems: New challenges in the llm era," in *Proc. 30th ACM SIGKDD Conf. Knowledge Discovery and Data Mining*, 2024, pp. 6437–6447.

[2] A. Défossez, L. Mazaré, M. Orsini, A. Royer, P. Pérez, H. Jégou, E. Grave, and N. Zeghidour, "Moshi: A speech-text foundation model for real-time dialogue," *arXiv preprint arXiv:2410.00037*, 2024.

[3] W. Chen, Z. Ma, R. Yan, Y. Liang, X. Li, R. Xu, Z. Niu, Y. Zhu, Y. Yang, and Z. Liu, "Slam-omni: Timbre-controllable voice interaction system with single-stage training," *arXiv preprint arXiv:2412.15649*, 2024.

[4] Q. Fang, S. Guo, Y. Zhou, Z. Ma, S. Zhang, and Y. Feng, "Llama-omni: Seamless speech interaction with large language models," *arXiv preprint arXiv:2409.06666*, 2024.

[5] J. C. Yang, D. Dalisan, M. Korecki, C. I. Hausladen, and D. Helbing, "Llm voting: Human choices and ai collective decision-making," in *Proc. AAAI/ACM Conf. AI, Ethics, Soc.*, 2024, vol. 7, pp. 1696–1708.

[6] X. Lin, W. Wang, Y. Li, S. Yang, F. Feng, Y. Wei, and T.-S. Chua, "Data-efficient fine-tuning for llm-based recommendation," in *Proc. 47th Int. ACM SIGIR Conf. Res. Develop. Inf. Retr.*, 2024, pp. 365–374.

[7] S. K. Sakib and A. B. Das, "Challenging fairness: A comprehensive exploration of bias in llm-based recommendations," in *Proc. IEEE Int. Conf. Big Data (BigData)*, 2024, pp. 1585–1592.

[8] I. O. Gallegos, R. A. Rossi, J. Barrow, M. M. Tanjim, S. Kim, F. Dernoncourt, T. Yu, R. Zhang, and N. K. Ahmed, "Bias and fairness in large language models: A survey," *Comput. Linguistics*, vol. 50, no. 3, pp. 1097–1179, 2024.

[9] H. Kotek, R. Dockum, and D. Sun, "Gender bias and stereotypes in large language models," in *Proc. ACM Collective Intelligence Conf.*, 2023, pp. 12–24.

[10] M. K. Ngueajio and G. Washington, "Hey asr system! why aren't you more inclusive? automatic speech recognition systems' bias and proposed bias mitigation techniques. a literature review," in *Int. Conf. Human-Computer Interaction*. Springer, 2022, pp. 421–440.

[11] C.-Y. Kuan and H.-Y. Lee, "Gender bias in instruction-guided speech synthesis models," *arXiv preprint arXiv:2502.05649*, 2025.

[12] M. Jahan, P. Mazumdar, T. Thebaud, M. Hasegawa-Johnson, J. Villalba, N. Dehak, and L. Moro-Velazquez, "Unveiling performance bias in asr systems: A study on gender, age, accent, and more," in *ICASSP 2025-2025 IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*. IEEE, 2025, pp. 1–5.

[13] Z. Fan, R. Chen, T. Hu, and Z. Liu, "Fairmt-bench: Benchmarking fairness for multi-turn dialogue in conversational llms," *arXiv preprint arXiv:2410.19317*, 2024.

[14] I.-E. Veliche, Z. Huangqun, V. A. Kochaniyan, F. Peng, O. Kalinli, and M. L. Seltzer, "Towards measuring fairness in speech recognition: Fairspeech dataset," *arXiv preprint arXiv:2408.12734*, 2024.

[15] J. Richter, Y.-C. Wu, S. Krenn, S. Welker, B. Lay, S. Watanabe, A. Richard, and T. Gerkmann, "Ears: An anechoic fullband speech dataset benchmarked for speech enhancement and dereverberation," *arXiv preprint arXiv:2406.06185*, 2024.

[16] Y.-C. Lin, W.-C. Chen, and H.-Y. Lee, "Spoken stereoset: on evaluating social bias toward speaker in speech large language models," in *2024 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2024, pp. 871–878.

[17] C. Sekkat, F. Leroy, S. Mdhaffar, B. P. Smith, Y. Estève, J. Dureau, and A. Coucke, "Sonos voice control bias assessment dataset: A methodology for demographic bias assessment in voice assistants," *arXiv preprint arXiv:2405.19342*, 2024.

[18] M. R. Costa-jussà, M. C. Meglioli, P. Andrews, D. Dale, P. Hansanti, E. Kalbassi, A. Mourachko, C. Ropers, and C. Wood, "Mutox: Universal multilingual audio-based toxicity dataset and zero-shot detector," *arXiv preprint arXiv:2401.05060*, 2024.

[19] J. Zhao, T. Wang, M. Yatskar, V. Ordonez, and K.-W. Chang, "Gender bias in coreference resolution: Evaluation and debiasing methods," *arXiv preprint arXiv:1804.06876*, 2018.

[20] M. Nadeem, A. Bethke, and S. Reddy, "Stereoset: Measuring stereotypical bias in pretrained language models," *arXiv preprint arXiv:2004.09456*, 2020.

[21] N. Nangia, C. Vania, R. Bhalerao, and S. R. Bowman, "Crows-pairs: A challenge dataset for measuring social biases in masked language models," *arXiv preprint arXiv:2010.00133*, 2020.

[22] Y. Li, H. Su, X. Shen, W. Li, Z. Cao, and S. Niu, "Dailydialog: A manually labelled multi-turn dialogue dataset," *arXiv preprint arXiv:1710.03957*, 2017.

[23] P. Budzianowski, T.-H. Wen, B.-H. Tseng, I. Casanueva, S. Ultes, O. Ramadan, and M. Gašić, "Multiwoz–a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling," *arXiv preprint arXiv:1810.00278*, 2018.

[24] K. Li, C. Shen, Y. Liu, J. Han, K. Zheng, X. Zou, Z. Wang, X. Du, S. Zhang, H. Luo, and et al., "Audiotrust: Benchmarking the multifaceted trustworthiness of audio large language models," *arXiv preprint arXiv:2505.16211*, 2025.

[25] J. Zhang, K. Bao, Y. Zhang, W. Wang, F. Feng, and X. He, "Is chatgpt fair for recommendation? evaluating fairness in large language model recommendation," in *Proc. 17th ACM Conf. Recommender Syst. (RecSys)*, 2023, pp. 993–999.

[26] Y. Chu, J. Xu, X. Zhou, Q. Yang, S. Zhang, Z. Yan, C. Zhou, and J. Zhou, "Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models," *arXiv preprint arXiv:2311.07919*, 2023.

[27] A. Zeng, Z. Du, M. Liu, K. Wang, S. Jiang, L. Zhao, Y. Dong, and J. Tang, "Glm-4-voice: Towards intelligent and human-like end-to-end spoken chatbot," *arXiv preprint arXiv:2412.02612*, 2024.

[28] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. J. Ostrow, A. Welihinda, A. Hayes, A. Radford, and et al., "Gpt-4o system card," *arXiv preprint arXiv:2410.21276*, 2024.

[29] G. Comanici, E. Bieber, M. Schaekermann, I. Pasupat, N. Sachdeva, I. Dhillon, M. Blistein, O. Ram, D. Zhang, E. Rosen, and et al., "Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities," *arXiv preprint arXiv:2507.06261*, 2025.

[30] S. Zhou, Y. Zhou, Y. He, X. Zhou, J. Wang, W. Deng, and J. Shu, "Indextts2: A breakthrough in emotionally expressive and duration-controlled auto-regressive zero-shot text-to-speech," *arXiv preprint arXiv:2506.21619*, 2025.

[31] Haohan Guo, Fenglong Xie, Frank K. Soong, Xixin Wu, and Helen Meng, "A multi-stage multi-codebook vq-vae approach to high-performance neural tts," *arXiv preprint arXiv:2209.10887*, 2022.

[32] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2023, pp. 28492–28518.

[33] C. Xu, J. Zhang, Z. Chen, C. Xie, M. Kang, Y. Potter, Z. Wang, Z. Yuan, A. Xiong, and Z. Xiong, "Mmdt: Decoding the trustworthiness and safety of multimodal foundation models," *arXiv preprint arXiv:2503.14827*, 2025.

[34] Y. Ge, X. Zhao, L. Yu, S. Paul, D. Hu, C.-C. Hsieh, and Y. Zhang, "Toward pareto efficient fairness-utility trade-off in recommendation through reinforcement learning," in *Proc. 15th ACM Int. Conf. Web Search Data Mining (WSDM)*, 2022, pp. 316–324.

[35] M. Tomlein, B. Pecher, J. Simko, I. Srba, R. Moro, E. Stefancova, M. Kompan, A. Hrckova, J. Podrouzek, and M. Bielikova, "An audit of misinformation filter bubbles on youtube: Bubble bursting and recent behavior changes," in *Proc. 15th ACM Conf. Recommender Syst. (RecSys)*, 2021, pp. 1–11.