



# Finding the Human Voice in AI: Insights on the Perception of AI-Voice Clones from Naturalness and Similarity Ratings

Linda Bakkouche<sup>†1</sup>, Charles McGhee<sup>1</sup>, Emily Lau<sup>1</sup>, Stephanie Cooper<sup>1</sup>, Xinbing Luo<sup>1</sup>, Madeleine Rees<sup>1</sup>, Kai Alter<sup>2</sup>, Brechtje Post<sup>1</sup>, Julia Schwarz<sup>†1</sup>

<sup>1</sup>University of Cambridge, United Kingdom

<sup>2</sup>Newcastle University, United Kingdom

lb983@cam.ac.uk<sup>†</sup>, julia.schwarz.ac@gmail.com<sup>†</sup>

## Abstract

AI-generated voice clones are important tools in language learning, audiobooks, and assistive technology, but often struggle to replicate key prosodic features such as dynamic  $F_0$  variation. The impact of these differences on speech perception remain underexplored. To address this, we conducted two behavioural tasks, evaluating listeners' ratings of naturalness and similarity for human speech, three AI voice clones (ElevenLabs, StyleTTS-2, XTTS-v2), and a 30%  $F_0$  variation condition. ElevenLabs was rated comparably to human speech, while StyleTTS-2 and XTTS-v2 received lower ratings. Reduced  $F_0$  variation also led to lower ratings, suggesting that prosody is key to perceived naturalness and similarity. Listener ratings were further influenced by speaker accent and sex, but not by AI tool experience. These findings suggest that prosodic features and speaker-specific characteristics could be drivers for the varying performance of AI-voice clones.

**Index Terms:** Speech Perception, Speech Synthesis, Human-Computer Interaction, Prosody

## 1. Introduction

### 1.1. Background

Artificial intelligence (AI) voice cloning technology has advanced significantly, creating synthetic voices that closely mimic human speech. This progress has driven diverse applications, including personalized language instruction with real-time feedback [1], audiobooks with efficient narration [2], and assistive technologies for individuals with speech impairments, such as the NAO Robot in speech therapy [3].

A major challenge in AI speech synthesis is reproducing prosody, i.e. variations in pitch, stress, and rhythm. Kulan-gareth and colleagues [4] demonstrated that listeners can reliably distinguish AI-generated voices from human ones when prosodic irregularities, such as unnatural speech shifts or disrupted timing, are combined with pauses that break the natural flow of speech. Fundamental frequency ( $F_0$ ) variation could have an additional effect on perceived naturalness in AI speech:  $F_0$  variation plays a pivotal role in speech, governing pitch modulation and shaping the expressiveness, emotional depth, and communicative intent of spoken language [5]. Indeed, insufficient  $F_0$  variation may result in mechanical-sounding speech, rated as less natural by listeners and affecting intelligibility [6, 7]. Despite efforts to improve prosodic control [7, 8], many current AI systems still produce speech with flat intonation and limited dynamic expression, indicating unresolved challenges in replicating natural prosody [5]. Previous studies have focused on individual systems without comparing performance across multiple AI models and without controlled prosodic manipulations. This narrow focus limits our understanding of how

different AI voice clones perform in comparison to one another and in comparison to fully natural speech with less engaging prosody. Additionally, it remains unclear whether judgments of speaker similarity and naturalness to human voices are influenced by factors such as speaker accent or sex, as well as listeners' AI experience, as frequent interaction with AI tools may shape sensitivity to subtle acoustic differences.

### 1.2. The current study

In the present study, we address how listeners perceive three state-of-the-art AI voice clones, focusing on the influence of key prosodic features on perceived naturalness and speaker similarity. Specifically, we aim to answer the following questions:

1. How do listeners perceive the naturalness and similarity of AI-generated voices compared to human speech?
2. How does variation in fundamental frequency ( $F_0$ ) contribute to the perception of naturalness and similarity?

We selected MUSHRA tests (Multiple Stimuli with Hidden Reference and Anchor) as the primary evaluation method, given their effectiveness in capturing fine-grained perceptual differences [9]. Participants were presented with speech samples and asked to rate each on a scale of 0 (low) to 100 (high). The inclusion of a Hidden Reference (here: Natural Human Speech) and an Anchor ensures that ratings reflect relative differences in perceived quality across conditions, and encourages participants to rank stimuli based on nuanced perceptual cues, which may be missed in simpler rating systems like Mean Opinion Scores (MOS).

We conducted two MUSHRA tasks: one assessing Naturalness (Exp. 1) and the other Similarity to a reference speaker (Exp. 2). Test stimuli included natural human speech (Hidden Reference), three state-of-the-art AI voice-clone models (ElevenLabs<sup>1</sup>, StyleTTS-2, XTTS-v2) [10, 11], a controlled prosodic manipulation retaining 30% of the  $F_0$  variation from the natural sample, and a noise-vocoded control (Anchor). Including three AI voice clones allowed us to compare performance between commercial (ElevenLabs) and open source models (Style TTS, XTTS-v2), and to distill broader trends in AI voice synthesis (Question 1). Including a controlled 30%  $F_0$  condition allowed us to explore the role of dynamic pitch modulation in a controlled manner (Question 2). We hypothesized that both AI voice clones and reduced  $F_0$  variability would receive lower ratings of Naturalness and Similarity compared to human speech, and that these effects might be further modulated by speaker accent and listeners' prior exposure to AI tools.

<sup>1</sup><https://elevenlabs.io/>

## 2. Methodology

### 2.1. Participants

Participants ( $N = 66$ ) were recruited and compensated through Prolific, and no participants took part in both tasks. Participants who failed to give the Anchor a score of 0 in more than 33% of trials, and participants who assigned 0 to conditions other than the Anchor in more than 33% were removed (Exp. 1:  $N = 2$ ; Exp. 2:  $N = 9$ ). The remaining participants were native speakers of British English without language, leaning or hearing impairments, retaining 28 subjects in the Naturalness task ( $f = 18$ ; *Mean Age* = 30.32, *SD* = 7.54), and 27 subjects in the Similarity task ( $f = 13$ ; *Mean Age* = 30.19, *SD* = 7.04).

### 2.2. Stimuli and conditions

The stimuli were three-second samples of speech, categorized into six conditions: **Natural human speech** produced by three different speakers (two female and one male); three AI-voice clones based on the natural voices (**ElevenLabs**, **Style TTS**, **XTTS-V2**); a controlled prosodic manipulation retaining 30%  $F_0$  variation of the natural voices; and a control condition that was designed to be the most unnatural / dissimilar condition (**Anchor**).

The AI-voice clone models were chosen to capture a range of commercial (ElevenLabs) and open source tools (StyleTTS-2, XTTS-v2) as well as different architectures (LLM-based = ElevenLabs, XTTS-v2; Diffusion-based = StyleTTS-2). All voice clones were high-scoring models on Huggingface's TTS Arena<sup>2</sup> at the time of testing (10/2024), which indicated that they were among the best performing TTS models available. The Natural Human Speech recordings consisted of 1-2 sentence long sections of novel children's stories, recorded by two Standard Southern British English (SSBE) speakers (one male, one female) and one female General American (GA) English speaker. Each voice cloning model was conditioned on a 15-second sample of these three speakers, the content of which was not included in any of the listening tests. The default settings for voice cloning from each model's API or Github page were used to synthesize all the test utterances, reflecting the likely use of these models for educational content generation. The 30%  $F_0$  variation condition was included to examine the role of dynamic prosody in the perception of naturalness and similarity. To create this condition, the pitch contours of natural speech were manipulated using the Praat plug-in Vocal Toolkit<sup>3</sup> to systematically reduce  $F_0$  variation to 30% of the original sample, without altering other prosodic features like duration or intensity. This condition served as an intermediary step, providing insights into how reduced prosodic variability affects perceived naturalness and similarity to human voices. For the Naturalness task (Exp. 1), the Anchor was produced with a diphone synthesizer<sup>4</sup> to create a prosodically unnatural rendition of the human voice, passed through a TanH distortion filter<sup>5</sup>, and finally low pass filtered with a cut-off at 500 Hz. Whilst there cannot be an explicit definition of unnatural speech, this method produced speech which was easily recognised as the anchor by the majority of participants (cf. *Section 2.1 Participants*). For the Similarity task (Exp. 2), the Anchor was a recording of the same

lexical content by one of the other speakers from our recordings.

### 2.3. Experiment procedure

The experiments were conducted online with Gorilla Experiment Builder [12] and each took 25-35 minutes to complete. Participants were administered, in order: a demographic questionnaire; one of the two listening tasks (Exp. 1: Naturalness task, Exp. 2: Similarity task); a questionnaire about their use of AI-voice tools; and a post-experiment feedback questionnaire, which asked participants to rate the difficulty of the task and their own concentration during the experiment.

After participants were familiarized with an example, they completed 21 trials in each task. Each trial presented the six conditions in randomized order. In Experiment 1 (Naturalness), participants were asked to judge how natural each sample sounded on a scale of 0 (completely unnatural) to 100 (completely natural). Participants were also given the instruction to give the Hidden Reference a maximum rating and the Anchor a minimum rating. In Experiment 2 (Similarity), participants were asked to judge how similar each sample sounded on a scale of 0 (definitely a different speaker) to 100 (definitely the same speaker), compared to a reference speaker of the same voice producing a different sentence.

### 2.4. Data analysis

Listeners' previous exposure to AI-generated voices may have had an impact on their perception of the voice clones presented in the current experiment. AI exposure was therefore collected in a questionnaire administered to all participants, and frequency of AI use was included in our analyses of Naturalness and Similarity Ratings. Data from the questionnaire on participants' self-reported use of AI voice assistants are summarised in *Section 3.1*.

Data from the Naturalness and Similarity tests were analyzed separately in R (Version 4.2.1) following the same procedures. First, data were averaged over subjects and analyzed with Friedman tests to determine significant differences between conditions, as recommended for MUSHRA tests [13]. Secondly, to test for potential differences between the three different speakers in the voice samples (Male British, Female British, Female American) and a potential influence of participants' self-reported AI usage (No Usage, Medium, High) on judgments, we removed the Anchor condition, scaled the data between 0-1, and applied an Arcsine transformation to reduce the mid-point skew. We then conducted a linear mixed-effects model with the *lme4* package [14] on the normally distributed data, including the three-way interaction Condition\*Speaker\*AI-Use as fixed effects and Subject as random intercept. We also included Speaker as random slope since we hypothesized that speaker differences could be a likely source of variability. The model was then optimized with AIC comparison. Variables were sum-to-zero contrast coded. Post-hoc pairwise comparisons are presented with Bonferroni corrections. All data and analysis scripts are available on OSF<sup>6</sup>.

## 3. Results

### 3.1. Questionnaire on AI use

Listeners' previous exposure to AI-generated voices was collected in a questionnaire administered to all participants, and

<sup>2</sup><https://huggingface.co/spaces/TTS-AGI/TTS-Arena>

<sup>3</sup><http://www.praatvocaltoolkit.com>

<sup>4</sup><https://github.com/tomasgoiba/diphone-synthesizer>

<sup>5</sup><https://github.com/iver56/audiomentations>

<sup>6</sup><https://doi.org/10.17605/OSF.IO/WF3HA>

frequency of AI use was included in our analyses of Naturalness and Similarity Ratings. The questionnaire was designed to assess participants' overall AI tool usage, aligning with our aim of examining whether prior AI experience influences voice perception. Participants were categorized into three groups based on their frequency of use: Low (never, monthly), Moderate (weekly, several times a week), and High (daily, several times a day). Interestingly, AI engagement patterns varied across these groups, with lower-frequency users displaying broader usage across different AI tools, while higher-frequency users exhibited more concentrated usage (Figure 1). A large proportion of low-frequency users reported engagement with a variety of tools, indicating a diverse, but less intensive, interaction with AI technologies. Moderate and high-frequency users showed distinct engagement patterns, with a relatively larger proportion using entertainment tools and virtual assistants compared to other applications (as well as communication tools in the high-frequency group). Notably, education-related AI tool usage was consistently the least utilized category across all levels of AI usage frequency. These patterns suggest that engagement with different AI tools may become more specialised with increasing overall use.

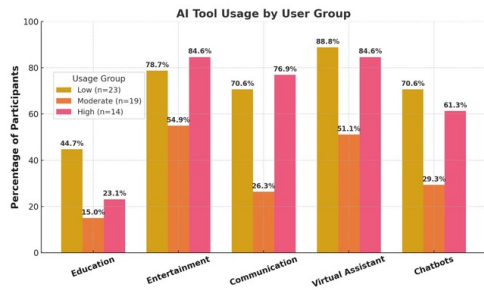


Figure 1: Participants % reporting engagement with different AI tools across different levels of overall AI-usage frequency.

### 3.2. Exp. 1: Naturalness task

In Experiment 1, listeners rated the naturalness of human and AI-voice cloned speech samples on a scale of 0-100. A Friedman test showed a significant effect of Condition on Naturalness ratings ( $\chi^2(5) = 131.08, p < .001$ ), with a large effect size (Kendall's  $W = .94$ ; Figure 2, left panel). Post-hoc pairwise Wilcoxon tests with Bonferroni correction revealed that all conditions differed significantly from one another except the Human voice and ElevenLabs conditions ( $W = 101, p = .312$ ), and the XTTS-V2 and 30%  $F_0$  conditions ( $W = 221, p = 1.000$ ). This suggests that ElevenLabs was the only voice-clone model that achieved human speech-like performance.

Next, a linear mixed model was built with the three-way interaction Condition\*Speaker\*AI-Use to investigate effects of Speaker Voice and AI Usage. AI-Use did not significantly contribute to the model and was removed. The resulting model showed significant effects of Condition ( $F(4, 378) = 254.93, p < .001, \eta^2 = .73$ ) and Speaker ( $F(2, 378) = 20.93, p < .001, \eta^2 = .10$ ), as well as a significant interaction between them ( $F(8, 378) = 5.55, p < .001, \eta^2 = .11$ ). Post-hoc comparisons of the interaction (with Bonferroni correction; emmeans package) showed two patterns of significant differences. In the Human natural, Human 30%  $F_0$ , and ElevenLabs conditions, Male British (MB) and Female British (FB) speakers were rated as significantly more natural than the Female American (FA)

speaker ( $ps < .05$ ), suggesting an effect of accent in these conditions. Secondly, in the 30%  $F_0$  condition, MB also significantly differed from FB, with MB being rated more natural compared to FB, indicating an additional effect of speaker differences, possibly driven by the speaker sex.

### 3.3. Exp. 2: Similarity task

In Experiment 2, listeners rated the similarity of human and AI-voice cloned speech samples compared to a reference speaker on a scale of 0-100. We also found a significant effect of Condition on Similarity ratings ( $\chi^2(5) = 117.05, p < .001$ ), with a large effect size (Kendall's  $W = .87$ ; Figure 2, right panel). Post-hoc comparisons with Bonferroni correction revealed that all conditions differed significantly from one another except the Natural Human voice and ElevenLabs conditions ( $W = 300, p = .097$ ), the StyleTTS and 30%  $F_0$  conditions ( $W = 185, p = 1.000$ ), and the XTTS-V2 and 30%  $F_0$  conditions ( $W = 71, p = .053$ ). This suggests that ElevenLabs was again the only voice clone that achieved human speech-like performance in terms of similarity to an original speaker.

A linear mixed model investigated potential effects of Speaker Voice and AI Usage. AI-Use did not significantly contribute to the model and was removed. The resulting model showed significant effects of Condition ( $F(4, 364) = 193.35, p < .001, \eta^2 = .68$ ), a non-significant effect of Speaker ( $F(2, 364) = 2.45, p = .087, \eta^2 = .01$ ), and a significant interaction between Condition and Speaker ( $F(8, 364) = 28.91, p < .001, \eta^2 = .39$ ). Post-hoc comparisons of the interaction showed that all contrasts with the Female American (FA) speaker were significant ( $ps < .01$ ). For the Natural Human, ElevenLabs, and 30%  $F_0$  samples, similarity ratings were higher on average for British speakers compared to the American speaker. Conversely, for the StyleTTS and XTTS-V2 samples, similarity ratings were lower for British speakers compared to the American speaker. In addition, only in the 30%  $F_0$  condition, British speakers (MB-FB) also differed significantly from one another ( $t(364) = 3.65, p < .001$ ), with MB being rated more similar to the reference speaker compared to FB, thus showing a similar pattern as in the Naturalness task. Listeners' accent experience, as well as differences in the underlying training data of the different AI-voice clone models, likely contributed to this pattern of results.

## 4. Discussion

Despite advances in AI voice cloning, replicating natural human speech remains challenging. However, research on listeners' perception of AI-voice clones remains sparse and has predominantly focused on assessing individual voice clone models in isolation [15, 16], limiting broader insights into AI perception. Our study addressed these gaps with two MUSHRA tasks in which we compared three state-of-the-art AI models (ElevenLabs, StyleTTS-2, XTTS-v2) to fully Natural Human Speech and to a less prosodically expressive, but also natural 30%  $F_0$  condition, to garner insights into the role of  $F_0$  variation.

Among the tested models, ElevenLabs was the only one rated comparably to human speech in both tasks, while StyleTTS-2 and XTTS-V2 were rated less natural and similar compared to the original speaker, consistent with studies showing AI struggles to replicate human voice variability [17]. Interestingly, reducing  $F_0$  variation to 30% also significantly affected ratings in both tasks, with comparatively similar ratings to the lower performing StyleTTS-2 and XTTS-V2 mod-

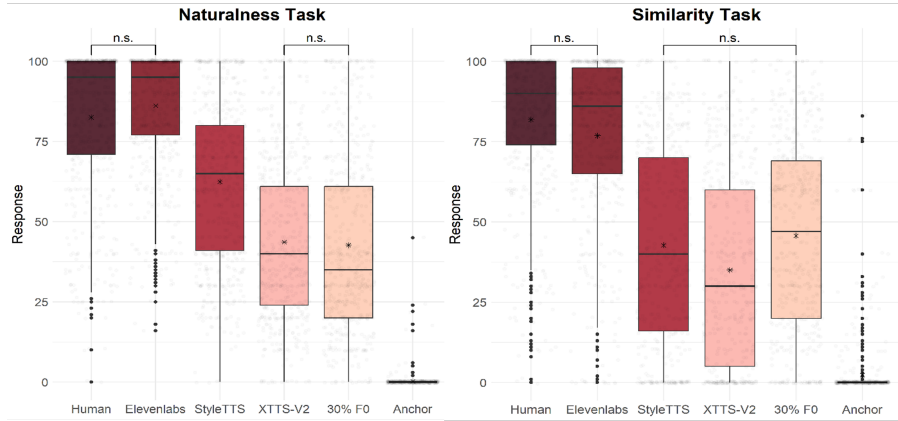


Figure 2: Listener Ratings of Naturalness (Exp. 1) and Similarity (Exp. 2).

els. While a causal claim cannot be made, this pattern of results suggests a potential role of natural  $F_0$  contours on the authenticity of AI speech. Indeed, research shows that listeners are highly sensitive to prosodic cues, effortlessly detecting even minor discrepancies that can undermine the perceived authenticity of speech [15]. Although beyond the scope of the current paper, the results prompt a detailed prosodic-phonetic analysis of AI-voice clones as the next step in this research, with likely differences in  $F_0$  contours.

Our results revealed further effects of speaker accent on listeners’ judgments: for Natural Human, ElevenLabs, and 30%  $F_0$  samples, British-accented speakers received higher similarity ratings than the American-accented speaker. Given that all our participants were British English speakers, this finding suggests that accent familiarity influenced judgments, consistent with judgments of natural speech [18]. Conversely, British speakers received lower ratings than the American speakers for the StyleTTS-2 and XTTS-V2 samples in the Similarity task. This was likely due to differences in training data: StyleTTS-2 is primarily trained on the LibriTTS dataset, which features predominantly American-accented speakers [10], and XTTS-v2 is trained on cross-lingual datasets with limited attention to British-accented prosody [11]. This may have affected the accurate replication of specific acoustic and prosodic features of British English, therefore affecting the Similarity to a British voice, but not the overall Naturalness. ElevenLabs’ flexibility and exposure to diverse datasets likely contributed to its superior handling of acoustic and prosodic nuances. Additionally, speaker sex appeared to influence listener judgments, with the male British speaker consistently receiving higher similarity ratings. This discrepancy may stem from the greater pitch variability in female voices compared to male ones [19], which can make precise voice matching more challenging. To improve generalization across different voices and enhance synthesis quality, models should be trained on properly curated, diverse datasets, as proposed by Luong and colleagues [20].

Our analysis of the AI usage questionnaire (Figure 1) revealed distinct engagement patterns: low-frequency AI users interacted with a broad array of tools, whereas high-frequency users focused on a narrower set of applications. However, this disparity in AI experience did not translate into differences in perceptual ratings – AI use frequency had no significant influence on listeners’ naturalness or similarity judgments. This finding suggests that even participants who regularly use AI voice tools were just as discerning in their evaluations as those

with little exposure, indicating that perceptual judgments were driven predominantly by the acoustic characteristics of the stimuli rather than by familiarity with the technology.

One limitation of our study is the short duration ( $\sim 12$ s) of the speech stimuli, which may not capture longer-range prosodic patterns. The relatively short sample length limits the analysis of broader prosodic dynamics such as phrasal intonation or boundary tones. Nonetheless, the comparison of prosodic implementation across synthesis systems remains valuable, and we plan to pursue extended phonetic analyses in future work.

## 5. Conclusions

This study revealed significant differences in human perception of current state-of-the-art voice clones: While ElevenLabs was rated comparably to human speech in both naturalness and similarity, StyleTTS-2 and XTTS-v2 received significantly lower ratings. Notably, a 30%  $F_0$  variation condition also led to low ratings, suggesting that dynamic prosody significantly impacts listeners’ judgments. Finally, speaker accent and sex had additional effects on listener ratings, likely related to raters’ linguistic backgrounds and limitations of AI speech technology. Taken together, the current results provide behavioural evidence that prosodic features and accent-specific characteristics affect human judgments of speech. This could be a potential driver for the varying performance of AI-voice clones. Although beyond the scope of the current paper, these results spark further research into the phonetic-prosodic characteristics of AI speech, such as detailed phonetic analyses on  $F_0$  variation, pitch range, speech rate, and intensity. Identifying which of these features contribute most to perceived differences between AI and human speech will help to explain why some AI-generated voices are perceived as less human.

## 6. Acknowledgements

The authors would like to thank the Cambridge Language Sciences Incubator Fund for funding this research, and the anonymous reviewers for their comments and suggestions.

## 7. References

- [1] B. Zou, H. Reinders, M. Thomas, and D. Barr, "Editorial: Using artificial intelligence technology for language learning," *Frontiers in Psychology*, vol. 14, pp. 1664–1078, 2023.
- [2] J. R. Rachels and A. J. Rockinson-Szapkiw, "The effects of a mobile gamification app on elementary students' Spanish achievement and self-efficacy," *Computer Assisted Language Learning*, vol. 31, no. 1-2, pp. 72–89, 2018.
- [3] G. Georgieva-Tsaneva, A. Andreeva, P. Tsvetkova, A. Lekova, M. Simonska, V. Stancheva-Popkostadinova, G. Dimitrov, K. Rasheva-Yordanova, and I. Kostadinova, "Exploring the potential of social robots for speech and language therapy: a review and analysis of interactive scenarios," *Machines*, vol. 11, no. 7, 2023.
- [4] N. V. Kulangareth, J. Kaufman, J. Oreskovic, and Y. Fossat, "Investigation of deepfake voice detection using speech pause patterns: Algorithm development and validation," *JMIR Biomedical Engineering*, vol. 9, 2024.
- [5] J. Kane, M. N. Johnstone, and P. Szewczyk, "Voice synthesis improvement by machine learning of natural prosody," *Sensors*, vol. 24, no. 5, 2024.
- [6] L. M. Slowiczek and H. C. Nusbaum, "Effects of speech rate and pitch contour on the perception of synthetic speech," *Human Factors*, vol. 27, no. 6, pp. 701–712, 1985.
- [7] J. Yamagishi and T. Kobayashi, "Average-voice-based speech synthesis using HSMM-based speaker adaptation and adaptive training," *IEICE TRANSACTIONS on Information and Systems*, vol. 90, no. 2, pp. 533–543, 2007.
- [8] J. M. Vojtech, J. P. Noordzij Jr, G. J. Cler, and C. E. Stepp, "The effects of modulating fundamental frequency and speech rate on the intelligibility, communication efficiency, and perceived naturalness of synthetic speech," *American Journal of Speech-Language Pathology*, vol. 28, no. 2S, pp. 875–886, 2019.
- [9] B. Series, "Method for the subjective assessment of intermediate quality level of audio systems," *International Telecommunication Union Radiocommunication Assembly*, vol. 2, 2014.
- [10] Y. A. Li, C. Han, V. Raghavan, G. Mischler, and N. Mesgarani, "StyleTTS 2: Towards human-level text-to-speech through style diffusion and adversarial training with large speech language models," *Advances in Neural Information Processing Systems*, vol. 36, pp. 19 594–19 621, 2023.
- [11] E. Casanova, K. Davis, E. Gölge, G. Gökmar, I. Gulea, L. Hart, A. Aljafari, J. Meyer, R. Morais, S. Olayemi, and J. Weber, "XTTS: a massively multilingual zero-shot text-to-speech model," in *Interspeech 2024*, 2024, pp. 4978–4982.
- [12] A. L. Anwyl-Irvine, J. Massonnié, A. Flitton, N. Kirkham, and J. K. Evershed, "Gorilla in our midst: An online behavioral experiment builder," *Behavior Research Methods*, vol. 52, pp. 388–407, 2020.
- [13] C. Mendonça and S. Delikaris-Manias, "Statistical tests with MUSHRA data," in *Audio Engineering Society Convention 144*. Audio Engineering Society, 2018.
- [14] D. Bates, M. Mächler, B. Bolker, and S. Walker, "Fitting linear mixed-effects models using lme4," *Journal of Statistical Software*, vol. 67, no. 1, pp. 1–48, 2015.
- [15] C. Roswadowitz, T. Kathiresan, E. Pellegrino, V. Dellwo, and S. Fröhholz, "Cortical-striatal brain network distinguishes deepfake from real speaker identity," *Communications Biology*, vol. 7, no. 1, 2024.
- [16] C. Edwards, A. Edwards, B. Stoll, X. Lin, and N. Massey, "Evaluations of an artificial intelligence instructor's voice: Social Identity Theory in human-robot interactions," *Computers in Human Behavior*, vol. 90, pp. 357–362, 2019.
- [17] Z. Hodari, C. Lai, and S. King, "Perception of prosodic variation for speech synthesis using an unsupervised discrete representation of F0," in *Proceedings of Speech Prosody 2020*, 2020, pp. 965–969.
- [18] S. Njie, N. Lavan, and C. McGettigan, "Talker and accent familiarity yield advantages for voice identity perception: A voice sorting study," *Mem Cogn*, vol. 51, no. 1, pp. 175–187, Jan. 2023.
- [19] A. P. Simpson, "Phonetic differences between male and female speech," *Language and Linguistics Compass*, vol. 3, no. 2, pp. 621–640, 2009.
- [20] H.-T. Luong, X. Wang, J. Yamagishi, and N. Nishizawa, "Training multi-speaker neural text-to-speech systems using speaker-imbalanced speech corpora," in *Interspeech 2019*. International Speech Communication Association, 2019, pp. 1303–1307.