

Inconsistent prosodies more severely impair speaker discrimination of Artificial-Intelligence-cloned than human talkers

Wenjun Chen¹, Xiaoming Jiang^{1,2}, Jingyi Ge¹, Shuwan Shan¹, Siyuan Zou¹, Yiyang Ding³

¹Institute of Linguistics, Shanghai International Studies University, Shanghai 201620, China

²Key Laboratory of Language Science and Multilingual Artificial Intelligence, Shanghai International Studies University, Shanghai 201620, China.

³School of Russian and Eurasian Studies, Shanghai International Studies University, Shanghai 201620, China

wenjun.chen@shisu.edu.cn, xiaoming.jiang@shisu.edu.cn

Abstract

AI algorithms designed to clone human speaker identity are reportedly capable of replicating human-specific vocal expression. However, whether listeners can identify a single speaker expressing varying emotive states as one individual remains unclear, particularly never in AI-to-AI pairings. This study asked thirty-six Chinese listeners to hear two consecutive clips and to judge whether identical speakers delivered pairs of Chinese sentences in human-only and AI-only scenarios, with the prosody of the first and second clips being incongruent or congruent. We found a decrease in the accuracy of identifying the same speaker under inconsistent prosody conditions compared to consistent ones, a trend evident in both human-to-human and AI-to-AI pairs. Meanwhile, correctly distinguishing between two speakers was more challenging than identifying a single speaker, with AI pairs reporting notably poorer performance than human-human pairs. When presented with pairs of speakers using consistent prosody, listeners demonstrated significantly slower reaction times when identifying two speakers. Our findings suggest that vocal prosodies can lead to within-speaker identity variation, in which listeners form average-based representations and still recognise the same speaker across prosodies. The findings about the reduced capability in speaker discrimination in AI voices provide supportive evidence for the ‘out-group homogeneity effect’ of AI voice perception.

Index Terms: voice cloning, speech synthesis, vocal confidence, speaker discrimination, voice identity

1. Introduction

Human voice conveys both short-term emotional states and long-term information like age and speaker identity, and such paralinguistic information can be decoded through computational models and human listeners [1, 2]. These paralinguistic cues can also be encoded by synthetic talkers, and listeners differently perceive the human and synthetic talkers in dimensions like truthfulness and powerfulness in terms of person perception as well as softness, squeakiness, slowness, nasality and liveliness in regard to speech qualities [3]. Recent studies have suggested a human emotional intimacy effect where audiobook users find human-narrated speech more enjoyable and can better attract their attention and arouse more positive emotional responses [4] – which could be attributed to synthetic algorithms’ inability to express human-specific speech prosodies. Another work has noted that voice cloning

services initially used for cloning speaker identity can capture and replicate the vocal confidence-related prosodic features (confident, doubtful, and neutral-intending), generating AI voices that share both speaker identity and vocal confidence with the original human talkers [5].

Emotional expressions and individual speaker identity seem to share a mutual articulation basis. The diverse anatomy and physiology of speech production mechanisms, such as the size and shape of the larynx and vocal tract, alongside control over articulatory muscles, fundamentally contribute to encoding a unique vocal identity by enabling person-specific distinctive sound production [6]. On the one hand, individuals with larger body sizes exhibit lower fundamental frequencies (F0) and longer vocal tract lengths (VTL), aspects that are often associated with perceptions of authority or dominance [7]. On the other hand, the modulation of VTL and F0 in response to different emotional states, such as confidence versus doubt, illustrates how speakers can intentionally or unconsciously manipulate their vocal qualities to convey complex psychological and emotional states [5, 8]. Meanwhile, the perception of VTL and F0 is relevant to decoding *who is talking* from the speech stream [9]. Hence, speakers’ social intention motivates how they sound, both in their emotional states and the impression of individual identity.

Yet, gaps related to listeners’ ability to decode identities across emotions remain. Lavan et al. (2019) shifted speakers’ VTL and F0 in the X-axis and Y-axis with *Praat* and exposed listeners to learning speaker identities away from the centre and tested if listeners recognised the never-heard voices in the centre as ‘old’. The experiment found listeners classified the never-heard voices as familiar, indicating the existence of an average-based representation of speaker identity [9]. Still, this study’s VTL and F0 manipulation was not associated with prosodic- or context-specific pragmatic intentions, such as fear. Xu and Armony (2021) designed a similar training-testing task that exposed listeners to three talkers’ identities in either neutral or fearful prosody in the training stage and tested if listeners could recognise the trained known speakers in one prosody still as ‘old’ when the speakers are expressing themselves in another prosody in the testing stage. They reported an accuracy lower than but close to the chance level [10]. This result seemed to suggest that if VTL and F0 modulations are associated with pragmatic intentions, i.e., fearful vs. neutral prosodies, listeners could not recognise the same talker’s identity across prosodies.

Against this background, we hypothesise that speech-prosody-led VTL and F0 modulation (not highly expressive

ones [11]) could shift the speaker's identity but within a range so that listeners still recognise talkers in different prosodies as the same talker. To test this, we employed the AX discrimination task rather than the previous training-testing approach, as the AX task is suited for identity comparison and recognition under challenging conditions. The paradigm involves comparing two sequentially played voices to determine whether they are expressed by one or two talkers, even in scenarios where speech is reversed or presented in an unfamiliar language [12]. Our study employed a design with 2 prosodies (confident vs. doubtful) * 2 sources (pre-clone human vs. post-clone AI) * 2 prosody pairs (inconsistent vs. consistent) * speaker pairs 2 (same vs. different) design. Listeners were tasked to decide whether a pair of voices was produced by the same speaker or not, ignoring the other differences between the sounds, such as prosodies. In our study, the 'consistent' vs. 'inconsistent' design is illustrated by pairing two voices, A and X, where A features characteristics like 'human speaker & confident prosody', and X can either share these characteristics (consistent) or differ, for example, 'human speaker & doubtful prosody' (inconsistent), thus manipulating prosody consistency.

2. Speaker discrimination study

2.1. Methods

2.1.1. Participants

Thirty-six native Mandarin Chinese speakers, university students (26 females/10 males; Mean \pm SD Age: 20.52 ± 2.62 years for females, 21.00 ± 2.72 years for males; Years of Education: 17.04 ± 2.13 for females, 17.55 ± 1.97 for males) without reported auditory or mental impairments participated in the perception experiment. Compensation was set at 50 RMB per hour. The Ethics Committee of the Institute of Linguistics, Shanghai International Studies University approved the study.

2.1.2. Material and paradigm

The auditory stimuli were selected from a validated audio corpus (24 speakers; working paper). The audio selected ensured that for the same sentence, for instance, *I can fill in the form for you* (in Chinese), the AI-cloned voice and the human-produced voice for both confident and doubtful prosodies were perceptually distinct in confidence level, with the confident condition scoring higher than the doubtful one. All 24 speakers were paired based on biological sex to form 12 contrasting pairs, also ensuring close similarity in speaker height within each pair.

The instruction was delivered in person. During the briefing, participants were instructed to ignore the linguistic content, prosodical differences, and the possible perceived unnaturalness in the speech pairs, only focusing on judging if one or two talkers were talking. Listeners were presented with a beep sound followed by two consecutive audio clips, namely Sound A and X, and were required to judge whether the same person spoke the two sentences by pressing keys (F or J). The assignment of F and J to responses was counterbalanced, as was the order of AX sentences across participants. Another manipulation is the *Prosody Consistency*: Consistent Prosody (CP) and Inconsistent Prosody (IP). The pairs were from the same source, either human-human or AI-AI. An additional manipulation is the sequence. For a speaker pair (e.g., Speaker 1 and 2), we counterbalanced which speaker's voice was presented first (A) and which was presented second (X). Each participant evaluated 192 pairs, encompassing 12 speaker pairs,

2 levels of prosody consistency (IP vs. CP), 2 sources (human-human or AI-AI), 2 types of talker numbers (one or two), and 2 arrangements of speaker sequences (AX and XA). The block orders were also balanced among participants. See Figure 1.

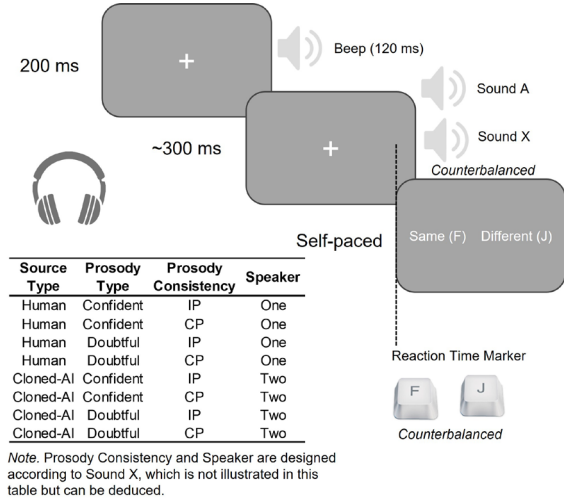


Figure 1: Design of the task.

2.1.3. Data analysis

The Effects of Talker and Prosody Consistency on Accuracy. With the *lme4* package (Version 1.1-35.1)[13] on RStudio version 2023.09.1 Build 494, our initial Mixed-effects logistic regression (MELR) model (Model 1-a) utilised the formula of *Response Accuracy ~ Prosody Consistency (CP and IP) * Talker (One and Two) + (1 | Participant)*, family = binomial. For 'Talker', the response accuracy corresponds to the correct answer for the pair: one talker (equals the 'same') and two talkers (equals the 'different'). The main effect and interaction results are documented in Table 1. Significant post hoc results are annotated in Figure 2 (contrasting one vs. two speakers in CP-CP or IP-IP). We further included speaker identity sources into the model (Model 1-B) with a formula: *Response Accuracy ~ Prosody Consistency (CP and IP) * Talker (One and Two) * Sources (human vs. AI) + (1 | Participant)*, family = binomial in a second fitting. Sources were added to explore the interaction between the AI-AI/human-human pairs and other factors. Post hoc results are reported in Table 2, annotated in Figure 2 (contrasting AI and human trials in sub-comparison).

The Effects of Talker and Prosody Consistency on Reaction Time (RT). We first fitted an LMER model (Model 2-A): *Reaction Time ~ Prosody Consistency (CP and IP) * Talker (One and Two) + (1 | Participant)*. Post hoc results are annotated in Figure 3 (contrasting one vs. two speakers in CP-CP or IP-IP). We fitted another model (Model 2-B) with a formula of *Reaction time ~ Prosody Consistency (CP and IP) * Talker (One and Two) * Sources (human vs. AI) + (1 | Participant)*. Sources are added to explore if they interact with other factors. Post hoc results are reported in Table 4, annotated in Figure 3 (contrasting AI and human trials in sub-comparison).

3. Results

3.1. Prosody Consistency and Talker Sources on Accuracy

For Model 1-A, prosody consistency and talker numbers influenced the likelihood of correctly identifying paired stimuli within human- or AI-only contexts (Table 1).

The IP level of the *Prosody Consistency* variable, relative to the CP baseline, revealed a substantial effect. An β of -2.64 implied lower log odds for the IP condition compared to CP. The odds ratio of .07 indicated that the chance of a correct answer under the IP was 7% of that under CP, suggesting that IP had a lower accuracy than CP.

In scenarios requiring participants to differentiate between voices from two distinct talkers, as opposed to a single speaker, the model suggested a notable decrease in the likelihood of a correct response. An β of -3.08 indicated a significant reduction in accuracy with two speakers. The corresponding Odds Ratio of .05 indicated that the probability of a correct response with two speakers was merely 5% of the likelihood with one speaker, highlighting the tendency of perceiving two speakers as one.

The interaction effect ($\beta = 3.08$) in the IP condition with two speakers significantly increased the odds of a correct response compared to one speaker in the CP baseline. However, this does not mean that the overall accuracy of IP-Two was necessarily higher than that of CP-One. This is because of the negative independent effect of ‘Two’ ($\beta = -3.08$), which suggested that discriminating between more speakers generally lower response accuracy. Thus, while the interaction suggested a relative improvement under IP with two speakers, the actual overall accuracy rate depended on the combined effects of all factors, not solely on the interaction.

Table 1: MELR Results for Model 1-a

Term ^a	β	SE	z	p ^b	OR ^c	95% CI
Intercept	4.33	.22	20	***	76.11	[50.85-119.65]
IP	-2.64	.22	-11.99	***	.07	[.04-.11]
Two	-3.08	.21	-14.42	***	.05	[.03-.07]
IP:Two	3.08	.23	13.22	***	21.75	[14.01-35.16]

^a IP (from IP & CP); Two (from One & Two Talkers)

^b Significance codes: $p < .001$ ***, $p < .01$ **, $p < 0.05$ *

^c The ‘Odds Ratio’ column is the exponentiated version of the estimates, which gives the change in odds for a one-unit increase in the predictor variable.

For Model 1-B (Table 2 and Figure 2), the newly added *Sources* (human vs. AI) had no main effect ($p=.22$), no interaction with *Prosody Consistency* ($p=.15$) and *Talker* ($p=.55$), and no three-level interaction with others ($p=.14$). Still, we observed that AI-AI and human-human pairs could influence listeners’ performance in detailed one vs. two talker-discrimination tasks through post hoc analysis. The results suggested that listeners’ performance of accurately discriminating between two talkers was higher in human-human pairs than in AI-AI pairs.

Table 2: Post Hoc Results for Model 1-B

Contrast ^a	β	SE	z	p ^b
One CP AI - One IP AI	2.36	.28	8.39	***
One CP AI - One CP H	-.52	.42	-1.23	.92
One IP AI - One IP H	.12	.13	.93	.98
Two CP AI - Two IP AI	-.42	.12	-3.43	*
Two CP AI - Two CP H	-.78	.12	-6.48	***
Two IP AI - Two IP H	-.85	.14	-6.09	***
One CP H - One IP H	3.00	.35	8.58	***
Two CP H - Two IP H	-.49	.15	-3.21	*

^a One CP AI: one AI speaker talking in consistent prosody in the pair. H for Human.

^b Significance codes: $p < .001$ ***, $p < .01$ **, $p < .05$ *

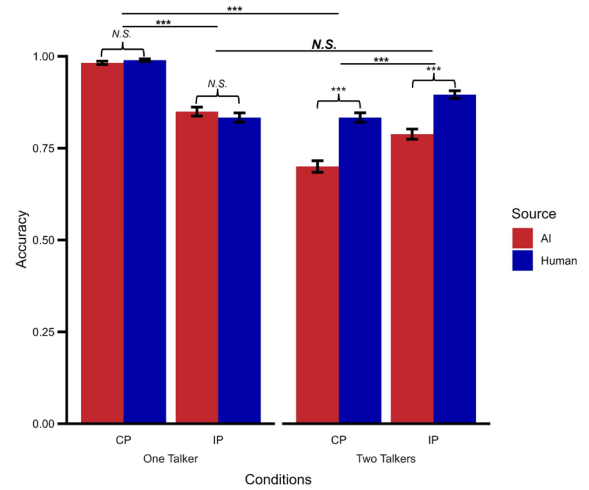


Figure 2: Talker numbers, prosody consistency, and speakers’ group identity influence accuracy.

3.2. Prosody Consistency and Talker Sources on RT

Analysis of Model 2-A revealed that the sources, whether AI or human, had no statistically significant main effect on participants’ reaction times ($F(1, 2280.4) = .93, p=.34$).

The talker numbers had a main effect, suggesting that reaction times varied depending on whether there was one or two talkers present ($F(1, 6776.1)=43.64, p < .001, \eta^2= 6.40e-03$). Post hoc contrast suggested that listeners reacted significantly slower when there were two talkers talking (One-Two: $\beta=-.23, SE=.03, z=-6.61, p<.0001$).

There was a significant interaction between the sound source and the number of talkers ($F(1,6776.1)=22.67, p < .001, \eta^2= 3.33e-03$), indicating that the effect of the number of talkers on reaction time was influenced by whether the sound source was AI or human. This indicated that listeners’ reaction times under conditions of one or two speakers were differentially affected by whether it was an AI-AI pair or a human-human pair. Additionally, following a post hoc analysis, we found listeners reacted significantly slower in the two-talkers condition only in IP rather than CP (CP One - CP One: $\beta=-.21, SE=.06, z=-3.74, p=.001$).

For Model 2-B, the post hoc analyses, as indicated by the annotated signs in Figure 3 and results in Table 3, demonstrated that no condition displayed a significant difference. We observed neither the main effect of the newly added human vs. AI sources ($p=.54$) nor its interaction with *Prosody Consistency* ($p=.78$) and *Talker* ($p=.19$), nor three-level interaction ($p=.99$). Further, post hoc analysis contrasting AI-AI and human-human pairs did not report any significance for CP-One ($p=1.0$), IP-One ($p=1.0$), CP-Two ($p=.96$), and IP-Two ($p=1.0$). Despite this, listeners were seemingly faster in AI pairs than human pairs only in one-speaker conditions, whereas they were slower in two-speaker conditions. However, this was supported only by visualisation ($p>.005$).

Table 3: Post Hoc Results for Model 2-B

Contrast ^a	β	SE	z	p
One CP AI - One IP AI	-.2	.08	-2.68	.13
One CP AI - One CP H	-.02	.07	-.22	1

One IP AI - One IP H	-.03	.07	-.48	1
Two CP AI - Two IP AI	.14	.08	1.81	.61
Two CP AI - Two CP H	.08	.07	1.1	.96
Two IP AI - Two IP H	.06	.07	.83	.99
One CP H - One IP H	-.22	.08	-2.92	.07
Two CP H - Two IP H	.12	.08	1.54	.79

^a One CP AI: one AI speaker talking in consistent prosody in the pair. H for Human.

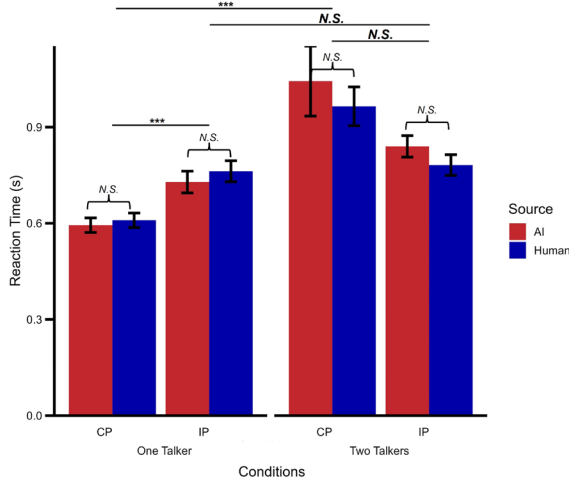


Figure 3: Talker numbers, prosody consistency, and speakers' group identity influence reaction time.

4. Discussion

We observed that listeners' performance in identifying whether two audio segments belong to one single speaker or rejecting them as two different people is influenced by prosody consistency and speaker group identity being AI or human.

Firstly, listeners' accuracy in identifying the same speaker as one talker significantly decreases with inconsistent prosody compared to consistent prosody. This pattern suggests that the identity discrimination was impaired due to prosody inconsistency. However, the accuracy of our AX study is above the chance level, which is distinct from identity recognition tasks which tasked participants to classify learned talkers as old or new across fearful vs. neutral prosodies [10] or group two characters (from the TV series *Breaking Bad*) with a highly expressive speech (e.g., shouting or strained voice) [11]. This observed improved integration performance might be related to the lower cognitive resources demanded in the AX discrimination task, which does not require listeners to carry forward speaker identity information for later recognition [12], unlike the training-testing task that demands listeners to become thoroughly familiar with individual speaker identities to adapt to the internal identity changes brought by audio VTL/F0 modulations [14]. Thus, the AX paradigm in the current study is suitable for directly exploring listeners' integration of speaker identity across speech prosody, similar to talker identification across languages [15]. Apart from the paradigm suitability, our study also adds to the human-like nature of AI-generated audio, as we found similar patterns in both AI and human trials - the prosody inconsistency impairs listeners' discrimination of talker identities. Hence, our study suggests that while inconsistent prosody can complicate speaker recognition, listeners can still effectively integrate one

talker in different prosodies. The similarity between AI-generated voices and human speakers in terms of identity perception suggests that these AI voices are crafted in ways that human listeners can easily interpret, potentially amplifying AI's utility in contexts requiring vocal interaction with humans [16].

Secondly, listeners were more likely to correctly identify two speakers in IP than in CP. This might reflect a tendency: when prosody is inconsistent, the perceived distance between speaker identities is greater, e.g., confident prosody has an averagely longer VTL and shorter F0 than doubtful prosody [5, 8]. Lavan et al. (2019) concluded that when stimuli were acoustically closer to the centre or average of a voice identity's representation, the accuracy in identifying these voices increased [9]. Our studies suggest that two talkers' VTL and F0 are more distinct when they express different emotive states [5], leading listeners to reject them as two speakers more easily.

Thirdly, when listeners should reject two speakers correctly as different individuals, performance in AI-AI conditions was significantly less accurate than in human-human conditions, in both IP and CP. This suggests that while being presented with two speakers, listeners are more inclined to consider AI speakers as single individuals, even in setups involving two speakers. This tendency might be related to the categorisation perception or the out-group homogeneity effect, where people perceive members of an out-group as more similar than members of their own group (in-group) [17]. Hence, less familiar or categorised groups (i.e., AI voices) are perceived as less distinct. Still, caution should be taken since omitting explicit mention of AI voices in the briefing before the experiment might affect the categorisation process, despite listeners later orally report their speculation about AI voices.

Finally, listeners need significantly more time to make a decision when they are presented with two speakers with consistent prosody, but this effect is not seen when prosody is inconsistent. In our current design, the reaction time was calculated from the end of the second audio segment, and this could lead to delayed detection of online speaker discrimination. After all, previous studies have suggested that speaker group identity can be differentiated as early as 100ms, as indicated by the N100 amplitude being sensitive to out-group attitudes [18]. We suppose the longer reaction time for telling two speakers apart can be related to the higher cognitive load required to process two different speaker identities. In the case of inconsistent prosody, the lack of significant reaction time differences between one and two talkers might be due to the inherently slower responses in this condition, possibly because processing prosody and distinguishing between speakers may occur in parallel rather than in an additive manner [18]. This hypothesis could be further supported by data that is sensitive to timing, specifically by using the beginning of the second audio segment as a reference point.

5. Acknowledgements

This work was supported by the Natural Science Foundation of China (Grant No. 31971037); the 'Shuguang Programme' supported by the Shanghai Education Development Foundation and Shanghai Municipal Education Committee (Grant No. 20SG31); the Natural Science Foundation of Shanghai (22ZR1460200); the Supervisor Guidance Programme of Shanghai International Studies University (2022113001); and the Major Programme of the National Social Science Foundation of China (Grant No. 18ZDA293).

6. References

- [1] B. Schuller and A. Batliner, *Computational paralinguistics: emotion, affect and personality in speech and language processing*. John Wiley & Sons, 2013.
- [2] M. Mileva and N. Lavan, "Trait impressions from voices are formed rapidly within 400 ms of exposure," *Journal of Experimental Psychology: General*, 2023.
- [3] J. W. Mullennix, S. E. Stern, S. J. Wilson, and C.-I. Dyson, "Social perception of male and female computer synthesised speech," *Computers in Human Behavior* vol. 19, no. 4, pp. 407-424, 2003.
- [4] E. Rodero and I. Lucas, "Synthetic versus human voices in audiobooks: The human emotional intimacy effect," *New Media & Society*, vol. 25, no. 7, pp. 1746-1764, 2023.
- [5] W. Chen and X. Jiang, "Voice-Cloning Artificial-Intelligence Speakers Can Also Mimic Human-Specific Vocal Expression," in *Preprints*, ed: Preprints, 2023.
- [6] R. J. Podesva and P. Callier, "Voice quality and identity," *Annual review of applied Linguistics*, vol. 35, pp. 173-194, 2015.
- [7] P. Sorokowski *et al.*, "Voice of authority: professionals lower their vocal frequencies when giving expert advice," *Journal of Nonverbal Behavior*, vol. 43, no. 2, pp. 257-269, 2019.
- [8] X. Jiang and M. D. Pell, "The sound of confidence and doubt," *Speech Communication*, vol. 88, pp. 106-126, 2017.
- [9] N. Lavan, S. Knight, and C. McGettigan, "Listeners form average-based representations of individual voice identities," *Nature Communications*, vol. 10, no. 1, pp. 1-9, 2019.
- [10] H. Xu and J. L. Armony, "Influence of emotional prosody, content, and repetition on memory recognition of speaker identity," *Quarterly Journal of Experimental Psychology*, vol. 74, no. 7, pp. 1185-1201, 2021.
- [11] N. Lavan, L. F. Burston, P. Ladwa, S. E. Merriman, S. Knight, and C. McGettigan, "Breaking voice identity perception: Expressive voices are more confusable for listeners," *Quarterly Journal of Experimental Psychology*, vol. 72, no. 9, pp. 2240-2248, 2019.
- [12] D. Fleming, B. L. Giordano, R. Caldara, and P. Belin, "A language-familiarity effect for speaker discrimination without comprehension," *Proceedings of the National Academy of Sciences*, vol. 111, no. 38, pp. 13795-13798, 2014.
- [13] D. Bates, "lme4: Linear mixed - effects models using Eigen and S4," *R package version*, vol. 1, p. 1, 2016.
- [14] N. Lavan, S. Knight, and C. McGettigan, "Listeners form average-based representations of individual voice identities," *Nature Communications*, vol. 10, no. 1, p. 2404, 2019/06/03 2019.
- [15] S. J. Winters, S. V. Levi, and D. B. Pisoni, "Identification and discrimination of bilingual talkers across languages," *The Journal of the Acoustical Society of America*, vol. 123, no. 6, pp. 4524-4538, 2008.
- [16] G. Laban, A. Kappas, V. Morrison, and E. S. Cross, "Building Long-Term Human-Robot Relationships: Examining Disclosure, Perception and Well-Being Across Time," *International Journal of Social Robotics*, 2023.
- [17] J. M. Ackerman *et al.*, "They all look the same to me (unless they're angry) from out-group homogeneity to out-group heterogeneity," *Psychological science*, vol. 17, no. 10, pp. 836-840, 2006.
- [18] X. Jiang, K. Gossack-Keenan, and M. D. Pell, "To believe or not to believe? How voice and accent information in speech alter listener impressions of trust," *Quarterly Journal of Experimental Psychology*, vol. 73, no. 1, pp. 55-79, 2020.