


AUTHOR QUERY FORM

	Journal: J. Acoust. Soc. Am.	Please provide your responses and any corrections by annotating this PDF and uploading it according to the instructions provided in the proof notification email.
	Article Number: JASA-12954	

Dear Author,

Below are the queries associated with your article; please answer all of these queries before sending the proof back to AIP.

Article checklist: In order to ensure greater accuracy, please check the following and make all necessary corrections before returning your proof.

1. Is the title of your article accurate and spelled correctly?
2. Please check affiliations including spelling, completeness, and correct linking to authors.
3. Did you remember to include acknowledgment of funding, if required, and is it accurate?

Location in article	Query / Remark: click on the Q link to navigate to the appropriate spot in the proof. There, insert your comments as a PDF annotation.
AQ1	Please check that the author names are in the proper order and spelled correctly. Also, please ensure that each author's given and surnames have been correctly identified (given names are highlighted in red and surnames appear in blue).
AQ2	Please verify the single-page article in Ref(s). Tajadura-Jimenez <i>et al.</i> (2017); otherwise, please provide the full page range.

Thank you for your assistance.



Beyond acoustics: Self-relevance as a key to voice naturalness (L)

AQ1

Ana P. Pinheiro^{a)}

Faculdade de Psicologia, CICPSI, Universidade de Lisboa, Lisboa, Portugal

ABSTRACT:

Synthetic voices can now achieve remarkable acoustic accuracy, yet often fail to sound “natural,” especially when designed to reproduce one’s own voice. Existing frameworks define naturalness along two dimensions: deviation from acoustic norms and human-likeness. Yet these dimensions overlook the self-voice, which can feel natural or unnatural for reasons beyond the signal itself. Here, self-relevance is proposed as a complementary dimension, capturing the subjective alignment between a voice and the listener’s self-representation. Evidence shows that self-relevance modulates perceived naturalness independently of acoustic match. A full understanding of voice naturalness, therefore, requires integrating physical speech properties with the listener’s self-representational framework.

© 2025 Acoustical Society of America. <https://doi.org/10.1121/10.0039927>

(Received 28 August 2025; revised 30 October 2025; accepted 3 November 2025; published online xx xx xxxx)

[Editor: Jody Kreiman]

Pages: 1–3

What makes a voice sound *natural*? For decades, acoustics has provided the primary explanation: naturalness is tied to how closely a signal conforms to expected patterns in the speech spectrum. Today, synthetic voices can achieve a near-perfect acoustic match to natural speech, yet still sound subtly “off” to listeners, often judged less favorably (Gessinger *et al.*, 2023; Herrmann, 2023) and, in particular, less socially appealing (Bruder *et al.*, 2025) than human voices. This paradox reveals that naturalness cannot be reduced to acoustics alone.

A recent framework (Nussbaum *et al.*, 2025) defines naturalness along two main acoustic-perceptual dimensions: *deviation-based naturalness* captures how much a voice diverges from typical acoustic patterns, whereas *human-likeness-based naturalness* reflects the extent to which a voice sounds recognizably human. In speech synthesis research, for example, naturalness is typically evaluated through perceptual methods such as the absolute category rating (ACR) and its outcome, the mean opinion score (MOS), which reflects the perceived closeness of synthetic to natural speech (Le Maguer *et al.*, 2024; Perrotin *et al.*, 2025). However, while the two dimensions explain why many artificial or degraded voices sound strange or unnatural, they overlook the self-voice, which can feel natural or unnatural for reasons extending beyond the signal itself. Here, I propose a complementary dimension of naturalness: self-relevance.

Whereas the *self-voice* refers to the acoustic stimulus associated with one’s own vocal identity, *self-relevance* captures the subjective alignment between that voice and the listener’s internal model of the self (Crow *et al.*, 2021). It is shaped by sensorimotor contingencies (Khalilian-Gourtani *et al.*, 2024), memory-based self-representations (Iannotti *et al.*, 2022), and multisensory integration

(Orepic *et al.*, 2023). Crucially, it is not reducible to *self-recognition*, i.e., the perceptual identification of a voice as one’s own (Candini *et al.*, 2014), or self-attribution, i.e., the inference that one is the source of the voice (Ohata *et al.*, 2022). Rather, it concerns whether a voice feels *natural for me*, even when acoustically atypical or externally generated. For instance, a hoarse version of one’s own voice during illness may sound more natural than a flawless synthetic clone.

Evidence supports treating self-relevance as separate from acoustic accuracy. Voices that deviate from typical acoustic features or patterns can still be judged as “natural” or “self-like” when they align with memory-based self-representations or multisensory predictions (Tajadura-Jiménez *et al.*, 2017), sometimes even altering self-perception (Tajadura-Jiménez *et al.*, 2017), mood (Aucouturier *et al.*, 2016), or social attitudes (Arakawa *et al.*, 2021). Conversely, acoustically unaltered voice feedback can feel “alien” or externally generated if it violates sensorimotor contingencies (Franken *et al.*, 2018) or lacks expected bone-conducted vibrotactile cues (Orepic *et al.*, 2023). Recent work (Rosi *et al.*, 2025b) further illustrates this complexity: participants sometimes rated cloned voices of others more favourably than their own, despite equivalent acoustic deviation. Such findings underscore that acoustics is necessary but not sufficient for judging a voice as *naturally “mine.”*

Neural evidence further supports the role of self-relevance. Like deviation and human-likeness, it modulates early auditory processing within ~200 ms after voice onset, consistent with rapid sensory gain for self-relevant input (Pinheiro *et al.*, 2016; Pinheiro *et al.*, 2023). These effects suggest that the brain treats the self-voice as a special category beyond simple acoustic analysis (e.g., Conde *et al.*, 2016, 2018). Its influence extends to later processing stages,

^{a)}Email: appinheiro@psicologia.ulisboa.pt

83 when voice input is evaluated against memory-based self-
84 representations (Iannotti *et al.*, 2022; Pinheiro *et al.*, 2023).
85 Therefore, self-relevance bridges (early) low-level auditory
86 and (later) higher-level cognitive analyses of voice informa-
87 tion. It may act as a contextual modulation of voice natural-
88 ness, explaining why some voice distortions feel authentic
89 or why even acoustically matched voice clones can fail to
90 sound “natural.”

91 Recognizing the limits of traditional evaluation criteria
92 (e.g., MOS; Le Maguer *et al.*, 2024), recent developments in
93 speech synthesis have reconceptualized naturalness as
94 *appropriateness*—the idea that synthetic speech should be
95 judged within its communicative context rather than in iso-
96 lation (Pandey *et al.*, 2025). This perspective aligns with
97 evidence showing that hearing is dynamically shaped by
98 perceptual, cognitive, social, and emotional context
99 (Kreiman, 2024). Within this framework, *self-voice evalua-*
100 *tion* represents a specific form of contextual appropriateness,
101 particularly relevant when synthesized voices aim to express
102 or preserve the user’s own vocal identity. Such cases include
103 personalized text-to-speech systems, clinical voice restora-
104 tion (e.g., for individuals with amyotrophic lateral sclerosis
105 or after laryngectomy), and personalized human-AI interfa-
106 ces (e.g., McGettigan *et al.*, 2024; Rosi *et al.*, 2025a).
107 Conversely, in contexts such as virtual assistants or public
108 service announcements—where the listener is not the voi-
109 ce’s owner—self-relevance becomes secondary to task-
110 related or social appropriateness.

111 Understanding when and how self-relevance contributes
112 to perceived naturalness has direct implications for both
113 clinical and technological domains. Altered self-voice feed-
114 back has been linked to auditory verbal hallucinations
115 (Pinheiro *et al.*, 2020), where acoustically realistic voices
116 are experienced as alien, not because they sound unnatural,
117 but because they lack the perceptual and neural signatures
118 of self-generation (e.g., Pinheiro *et al.*, 2020). In technologi-
119 cal contexts, while users often express positive attitudes
120 toward digital voice transformations (Gueraouou *et al.*,
121 2024), these technologies can have unintended psychologi-
122 cal consequences, revealing the limits of acoustics as a pre-
123 dictor of perceived naturalness. For example, voice cloning
124 may destabilize a speaker’s own sense of vocal identity,
125 consistent with evidence linking how we sound to who we
126 believe ourselves to be (Stern *et al.*, 2021).

127 To conclude, acoustics is essential but not sufficient for
128 understanding voice naturalness. Voice perception—including
129 that of synthetic voices—depends not only on the physi-
130 cal signal but also on cognitive and contextual factors that
131 shape its interpretation (Kreiman, 2024). Naturalness,
132 though still a loosely defined, multifaceted perceptual con-
133 struct (Pandey *et al.*, 2025), must therefore be understood as
134 emerging from the interaction between the acoustic proper-
135 ties of speech and the listener’s representational framework.
136 Introducing self-relevance as a third dimension highlights
137 the perceptual uniqueness of the self-voice and opens a
138 novel research agenda at the intersection of acoustics, self-
139 representation, and artificial intelligence—one that can

inform both scientific understanding and responsible techno- 140
logical design. For speech synthesis, this means expanding 141
beyond signal fidelity and human-likeness toward perceptual 142
congruence with the listener’s internal model of the self, 143
while developing evaluation methods that explicitly capture 144
this self-relevant dimension. In a world where the line 145
between genuine and synthetic speech grows even thinner, 146
accounting for self-relevance will be essential for under- 147
standing and preserving what makes a voice truly natural. 148

149 **ACKNOWLEDGMENTS**

This work was supported by Fundação para a Ciência e 150
Tecnologia and BIAL Foundation (Grant Nos. 2023.00041. 151
RESTART and BIAL 146/2020). The author would like to 152
thank the VoicES Lab members (www.voicesneurolab.com) 153
for the fruitful discussions about voice naturalness and its 154
underlying neural and functional mechanisms. 155
156

157 **AUTHOR DECLARATIONS**

158 **Conflict of Interest**

The author has no conflicts to disclose. 159

160 **DATA AVAILABILITY**

Data sharing is not applicable to this article as no new 161
data were created or analyzed in this study. 162
163

164 Arakawa, R., Kashino, Z., Takamichi, S., Verhulst, A., and Inami, M.
165 (2021). “Digital speech makeup: Voice conversion based altered auditory
166 feedback for transforming self-representation,” in *Proceedings of the*
167 *2021 International Conference on Multimodal Interact.*, pp. 159–167.
168 Aucouturier, J.-J., Johansson, P., Hall, L., Segnini, R., Mercadié, L., and
169 Watanabe, K. (2016). “Covert digital manipulation of vocal emotion alter
170 speakers’ emotional states in a congruent direction,” *PNAS* **113**, 948–
171 953.
172 Bruder, C., Breda, P., and Larrouy-Maestri, P. (2025). “Attractive synthetic
173 voices,” *Comput. Hum. Behav. Artif. Hum.* **6**, 100211.
174 Candini, M., Zamagni, E., Nuzzo, A., Ruotolo, F., Iachini, T., and
175 Frassinetti, F. (2014). “Who is speaking? Implicit and explicit self and
176 other voice recognition,” *Brain Cogn.* **92**, 112–117.
177 Conde, T., Gonçalves, Ó. F., and Pinheiro, A. P. (2016). “The effects of
178 stimulus complexity on the preattentive processing of self-generated
179 and nonself voices: An ERP study,” *Cogn. Affect. Behav. Neurosci.* **16**,
180 106–123.
181 Conde, T., Gonçalves, Ó. F., and Pinheiro, A. P. (2018). “Stimulus com-
182 plexity matters when you hear your own voice: Attention effects on self-
183 generated voice processing,” *Int. J. Psychophysiol.* **133**, 66–78.
184 Crow, K. M., van Mersbergen, M., and Payne, A. E. (2021). “Vocal congru-
185 ence: The voice and the self measured by interoceptive awareness,”
186 *J. Voice* **35**, 324.e15–324.e28.
187 Franken, M. K., Acheson, D. J., McQueen, J. M., Hagoort, P., and Eisner,
188 F. (2018). “Opposing and following responses in sensorimotor speech
189 control: Why responses go both ways,” *Psychon. Bull. Rev.* **25**, 1458–
190 1467.
191 Gessinger, I., Cohn, M., Cowan, B. R., Zellou, G., and Möbius, B. (2023).
192 “Cross-linguistic emotion perception in human and TTS voices,” in
193 *Proceedings Interspeech 2023*, pp. 5222–5226.
194 Gueraouou, N., Vaiva, G., and Aucouturier, J.-J. (2024). “Social affective
195 inferences in the era of AI-filters: Towards the Bayesian reshaping of
196 human sociality?,” HAL-04801092 (HAL, Lyon, France).
197 Herrmann, B. (2023). “The perception of artificial-intelligence (AI) based
198 synthesized speech in younger and older adults,” *Int. J. Speech Technol.*
199 **26**, 395–415.

- 200 Iannotti, G. R., Orepic, P., Brunet, D., Koenig, T., Alcoba-Banqueri, S.,
201 Garin, D. F. A., Schaller, K., Blanke, O., and Michel, C. M. (2022). "EEG
202 spatiotemporal patterns underlying self-other voice discrimination,"
203 *Cereb. Cortex* **32**, 1978–1992.
- 204 Khalilian-Gourtani, A., Wang, R., Chen, X., Yu, L., Dugan, P., Friedman,
205 D., Doyle, W., Devinsky, O., Wang, Y., and Flinker, A. (2024). "A corol-
206 lary discharge circuit in human speech," *Proc. Natl. Acad. Sci. U.S.A.*
207 **121**, e2404121121.
- 208 Kreiman, J. (2024). "Information conveyed by voice quality," *J. Acoust.*
209 *Soc. Am.* **155**, 1264–1271.
- 210 Le Maguer, S., King, S., and Harte, N. (2024). "The limits of the Mean
211 Opinion Score for speech synthesis evaluation," *Comput. Speech Lang.*
212 **84**, 101577.
- 213 McGettigan, C., Bloch, S., Rosi, V., Dinkar, T., Lavan, N., Bowles, C., and
214 Reus, J. (2024). "Voice cloning: Psychological and ethical implications of
215 intentionally synthesising familiar voice identities".
- 216 Nussbaum, C., Frühholz, S., and Schweinberger, S. R. (2025).
217 "Understanding voice naturalness," *Trends Cogn. Sci.* **29**, 467–480.
- 218 Ohata, R., Asai, T., Imaizumi, S., and Imamizu, H. (2022). "I hear my
219 voice; therefore I spoke: The sense of agency over speech is enhanced by
220 hearing one's own voice," *Psychol. Sci.* **33**, 1226–1239.
- 221 Orepic, P., Kannape, O. A., Faivre, N., and Blanke, O. (2023). "Bone con-
222 duction facilitates self-other voice discrimination," *R Soc. Open Sci.* **10**,
223 221561.
- 224 Pandey, A., Le Maguer, S., and Harte, N. (2025). "What is naturalness?," in
225 *Proceedings of the 13th edition of the Speech Synthesis Workshop*, pp.
226 215–221.
- Perrotin, O., Stephenson, B., Gerber, S., Bailly, G., and King, S. (2025).
227 "Refining the evaluation of speech synthesis: A summary of the Blizzard
228 Challenge 2023," *Comput. Speech Lang.* **90**, 101747. 229
- Pinheiro, A. P., Rezaei, N., Nestor, P. G., Rauber, A., Spencer, K. M., and
230 Niznikiewicz, M. (2016). "Did you or I say pretty, rude or brief? An ERP
231 study of the effects of speaker's identity on emotional word processing,"
232 *Brain Lang.* **153-154**, 38–49. 233
- Pinheiro, A. P., Sarzedas, J., Roberto, M. S., and Kotz, S. A. (2023).
234 "Attention and emotion shape self-voice prioritization in speech pro-
235 cessing," *Cortex* **158**, 83–95. 236
- Pinheiro, A. P., Schwartze, M., Amorim, M., Coentre, R., Levy, P., and Kotz,
237 S. A. (2020). "Changes in motor preparation affect the sensory consequences
238 of voice production in voice hearers," *Neuropsychologia* **146**, 107531. 239
- Rosi, V., Payne, B., and McGettigan, C. (2025a). "Effects of self-similarity
240 and self-generation on the perceptual prioritization of voices," *J. Exp.*
241 *Psychol. Hum. Percept. Perform.* **51**, 996–1007. 242
- Rosi, V., Soopramanien, E., and McGettigan, C. (2025b). "Perception and
243 social evaluation of cloned and recorded voices: Effects of familiarity and
244 self-relevance," *Comput. Hum. Behav. Artif. Hum.* **4**, 100143. 245
- Stern, J., Schild, C., Jones, B. C., DeBruine, L. M., Hahn, A., Puts, D. A.,
246 Zettler, I., Kordsmeyer, T., Feinberg, D., Zamfir, D., Penke, L., and
247 Arslan, R. C. (2021). "Do voices carry valid information about a speaker's
248 personality?," *J. Res. Personal.* **92**, 104092. 249
- Tajadura-Jiménez, A., Banakou, D., Bianchi-Berthouze, N., and Slater, M.
250 (2017). "Embodiment in a child-like talking virtual body influences object
251 size perception, self-identification, and subsequent real speaking," *Sci.*
252 *Rep.* **7**, 9637. 253