

# Trends in Cognitive Sciences

## Understanding Voice Naturalness

--Manuscript Draft--

<b>Manuscript Number:</b>	TICS-D-24-00198R2
<b>Article Type:</b>	Review
<b>Keywords:</b>	Naturalness; Human-likeness; Voice perception; Authenticity; Voice synthesis
<b>Corresponding Author:</b>	Christine Nussbaum Friedrich Schiller University Jena Jena, GERMANY
<b>First Author:</b>	Christine Nussbaum
<b>Order of Authors:</b>	Christine Nussbaum  Sascha Fröhholz  Stefan R. Schweinberger
<b>Abstract:</b>	Perceived naturalness of a voice is a prominent property emerging from vocal sounds, which affects our interaction with both human and artificial agents. Despite its importance, a systematic understanding of voice naturalness is elusive. This is due to (a) conceptual underspecification, (b) heterogeneous operationalization, (c) lack of exchange between research on human and synthetic voices, and (d) insufficient anchoring in voice perception theory. This review reflects on current insights into voice naturalness by pooling evidence from a wider interdisciplinary literature. Against that backdrop, it offers a concise definition of naturalness and proposes a conceptual framework rooted both in empirical findings and theoretical models. Finally, it identifies gaps in current understanding of voice naturalness and sketches perspectives for empirical progress.

Re: Your letter from January 09, 2025, per e-mail, regarding Manuscript ID TICS-D-24-00198R1 "Understanding Voice Naturalness"

Dear Lindsey,

we would like to thank you and the reviewers again for your helpful feedback on the above manuscript, and for inviting us to submit a revised version of our Review Paper "Understanding Voice Naturalness".

We included a response letter which gives details on how we addressed each individual point by the reviewers and how we incorporated your suggestions on the manuscript.

We are currently preparing the Supplemental Materials in accordance with the TiCS guidelines (see also my E-Mail from 15.01.2025). For now, we have not yet removed the references to the OSF repository, but we marked alongside how we would link the supplemental materials we are about to prepare.

We hope you agree that the revised version is now suitable for publication in *TiCS* and we look forward to your response.

With kind regards,

Christine Nussbaum (on behalf of all authors).

Comments by the editor:

1. **One big picture note: generally, Reviews should avoid using first-person language.**  
Although this is not a hard rule, the idea is that Reviews should not be centering the perspective of the authors, as this would be more suitable for an Opinion. Please try to revise the highlights and the abstract to place less focus on your perspective. This will give a stronger impression of a Review. For example, “We show that current voice naturalness research is situated within different research domains that resemble echo chambers within science – they neither cross-refer to one another nor to current voice perception theory” --> “Current voice naturalness research is situated within different research domains that resemble echo chambers within science. They neither cross-refer to one another nor to current voice perception theory”. You need not contort the language awkwardly to avoid the first person, but please tilt the language away from this, throughout the piece. I’ve edited the introduction section to provide an example.

Response:

Thank you for your edits and comments in the manuscript. We incorporated most of your suggested changes. In what follows, we will elaborate on a few ones where we slightly deviated from your recommendations:

Regarding first-person language: we revised the manuscript and the highlights accordingly. There are, however, a few exceptions where we kept the first-person expression. Specifically, we kept it in the section on the definitions of naturalness (e.g. “we propose a taxonomy”) because we feel that this reflects an original position and contribution of the authors and we therefore wanted to frame it as such. We hope, you agree with us on this point.

Regarding your question about the purpose of our argumentation based on the ChatGPT output (page 10): We have now rephrased that part to clarify the point we wished to make here, which now reads as follows: *“At first sight, the concepts of authenticity and naturalness appear highly similar. In fact, when ChatGPT was prompted for synonyms of naturalness, authenticity was its first reply (Figure 1B), which may suggest that in openly accessible online sources, these two terms are indeed frequently occurring in an interchangeable manner. Accordingly, it might be argued that authenticity is just a special form of deviation-based naturalness, with a more specific reference.”* (page 10)

On page 12, you asked whether the preliminary evidence about authenticity conflicted with the predictions made in the previous paragraphs: Yes, that is correct, and we appreciate that we should better highlight this point. We now amended: *“However, comparably early effects also have been....”* (page 12/13)

Figure 3, in the caption, you recommended replacing “enacted” with “performed”: Here, we prefer to keep the original wording, because “enacted” is an established word in the literature on vocal emotion perception, so readers may be more familiar with it.

Regarding the references: Thank you for pointing out that some information was missing. We rechecked all references and made some small additions. The literature we cite is very heterogeneous with regard to format due to the interdisciplinary nature of the topic, including books and conference contributions. Therefore, they sometimes deviate from classical literature about cognitive neuroscience (e.g. sometimes they have no page numbers). We compiled the citation information to the best of our knowledge, but it may be that some references may still need adjustment upon final review.

Regarding the supplemental materials: as of now, we have not yet removed the references to the OSF repository from the manuscript. However, we have marked alongside how we would link the supplemental materials we are about to prepare.

Comments by Reviewer 1:

**Thank you to the authors for their thorough work on revising the paper in response to my previous comments. On this reading I found the paper to be more coherent and impactful. It will be a valuable and timely contribution to the literature, and I look forward to seeing it published!**

Response: Thank you very much for this positive evaluation and your valuable feedback!

I just have a few very, very minor comments on wording:

2. **Page 6 line 20: "In principle, such empirical heterogeneity can be a powerful source of insight". I understood what is meant here, but somehow it didn't flow for me as it's immediately followed by a sentence suggesting that the impact of different scales might be minimal. So perhaps it requires a little more qualification about how heterogeneity might add to understanding rather than limiting it (or having no impact at all).**

Response: This is a valid point. We rephrased this part as follows:

*"In principle, such empirical heterogeneity can be a powerful source of insight, potentially revealing the degree to which methodological aspects affect results. For example, there is recent evidence from face perception that differences in rating scales may not have a large impact on outcome [66]. However, it cannot be concluded that this generalizes to naturalness ratings, and the insufficient report of empirical details impedes a meaningful comparison of findings."* (page 6)

3. **Page 6, line 35: "Finally, few studies only..." --> "Finally, only few studies.."**

Response: Changed as suggested.

4. **Page 6, line 48: "However, while the scientific findings are well-received within each domain" --> I'm not totally clear on the point that needs to be made here with the term "well-received". I usually take this term to mean that people like/appreciate something, but I'm not sure that can be said here without further qualification about something like the quality of the journals / citation rates in which the work appears. Is it that the work on naturalness is well-cited and/or theoretically impactful within individual research domains? Or more simply that it is appropriate to domain-specific audiences but doesn't translate across disciplinary boundaries?**

Response: Good point. Our point here specifically targets referencing and cross-citation, which seems to happen within domains but rarely across the disciplinary boundaries. To make this clearer, we reworded the sentence into:

*"However, while the scientific findings are acknowledged and referenced within each domain, these domains are poorly interconnected."* (page 6)

5. **Figure 2 caption: "...marked by the human voice border" --> I somehow just couldn't get this final sentence to scan easily. Are the samples defining the border, or rather (as I suspect) the human's perceptual evaluation of the samples? Would it be better to say something like "Human-likeness can be assessed from audio samples by judging whether**

**or not they lie within the perimeter of an acceptable human voice space (or "human voice border").**

Response: Changed as suggested.

- 6. Figure 3 caption: A suggestion for the penultimate sentence: "... or it could be real or fake in relation to person-related identity information."**

Response: Changed as suggested.

Comments by Reviewer 3:

**The authors have addressed most of my comments - I'd like to thank them for the thorough revision. There are just a few, minor issues that I'd like to push back on - just for the authors' consideration.**

Response: Thank you for this positive evaluation and your helpful feedback!

- 7. Spontaneous evaluation of voice naturalness - in my previous review I mentioned that I'm not aware of much evidence that listeners spontaneously evaluate voices on naturalness as some of the work on this topic does not seem to include judgements of naturalness when listeners are asked to freely describe their first impression of a particular voice. The authors offer a paper by Kuhne et al. (2020) in their response that does focus on naturalness, but it doesn't exactly address my point. In Kuhne et al. perceivers are explicitly asked to rate voices for naturalness or for how human they sound but they do not show any evidence that perceivers would evaluate voices on a dimension related to naturalness unprompted.**

Response: Indeed, the evaluations in Kühne et al. (2020) were promoted specifically for naturalness. Evidence for unprompted evaluations of naturalness is very sparse and presumably depends on the stimulus material used. In a dataset where healthy human, pathological human and synthetic voices are mixed, spontaneous evaluations related to naturalness may be more likely. Nevertheless, in the work of Lavan (2023), where the dataset was comprised of healthy human speakers only, some evaluations occurred that are related to naturalness, such as "robotic", "monotone", "speech impediment", or "muffled". Therefore, although more data would be desirable, there is evidence for unprompted naturalness evaluations. Lavan (2023) is referenced in the respective sentence in the manuscript.

- 8. Face/voice impression models - when discussing the different labels used for the main dimensions of face/voice person perception, the authors distinguish between warmth and competence (Fiske et al., 2007) and the trustworthiness-dominance model of Oosterhof and Todorov (2008). I think it's important to acknowledge that these models are based on different types of cues - while the latter model is based on face-based impressions, the former one is a more general social cognition model, mostly based on text cues. It, therefore, might be factually inaccurate to suggest that the warmth-competence model has been proposed for the underlying structure of face or voice impressions.**

Response: This is a good point. We have now clarified that these models in their original evolution were based on different types of cues. At the same time, it is important to note that there is now substantial work to integrate the warmth-competence model with the trustworthiness-dominance model (e.g., Sutherland et al., 2016, Cognition), validating our argument. In the paper, we now clarify as follows.

*"Poor interconnectivity is not unique to naturalness but can affect many other research domains within person perception. Consider fields with different research traditions, such as impression formation according to social psychological models of inter-group perception versus face/voice perception models. These models were developed for different types of perceptual cues, and different two-factor models with different labels have been proposed in both cases (e.g., warmth vs. competence, e.g. [70]; or trustworthiness vs. dominance, e.g. [71]). More recently though, these fields arguably benefited from interconnectivity, with substantial research to link these distinct clusters and uncover both these specific taxonomies and their empirical relationships [72,73]."* (page 6)

- 9. Deviation-based vs human-likeness-based naturalness - I still struggle a little to differentiate between these two types as one does seem like a version of the other. It is not quite clear to me why the human-likeness-based naturalness has the additional assumption of the existence of non-human voices and that's not true for deviation-based naturalness. Why couldn't the deviation-based naturalness cross the boundaries between human and artificial voices? Perhaps some further elaboration on this point would help clarify the taxonomy proposed by the authors.**

Response: This is a valid point. The key point here is that the additional assumption of the non-human voice is obligatory for the definition of human-likeness-based naturalness, whereas for the deviation-based naturalness it is not (although it is possible). This may seem as a small detail, but in our view, this has important practical implications and ties back to the intuitive understanding of readers regarding the distinction between pathological vs. synthetic voices. We tried to make this more specific by rephrasing the wording:

*"Compared to the deviation-based definition, the concept of human-likeness-based naturalness requires an additional obligatory assumption: the existence of a non-human voice space. This highlights the notion of a categorical boundary to human voices, although the transition between categories can be continuous. In other words, a definition of human-likeness is only meaningful if we assume that voices can be non-human in principle. Although deviation-based naturalness may, in certain cases, cross the boundary to the non-human voice space, this boundary is not essential for its definition. Apart from this critical distinction, however, human-likeness-based naturalness may represent a special case of deviation-based naturalness: the reference is a human voice (or listeners' representation of a human voice), and the deviation is assessed along the human/non-human spectrum. The above considerations suggest that the human-likeness-based conceptualization is particularly well-suited for research into synthetic voices."* (page8)

- 10. p.6, just before reference 67 - it is not clear what type of reliability the authors are referring to here**

Response: Good point. We specified this as "*interrater reliability*" now.

- 11. p. 6, last row - the authors refer to some highly divergent research traditions, but it is not clear what traditions they are referring to. It seems like they are referring to face vs voice impressions which I wouldn't necessarily call highly divergent, given that a lot of the voice impressions literature is strongly rooted and based on research on faces.**

Response: We have now reworded this, and replaced "highly divergent research traditions" with "different research traditions". Please also see our response to your point #8. Rephrasing this paragraph according to your suggestions also clarifies which different research traditions we refer to. Overall, we appreciate your careful and critical reading; we believe your comments have helped us clarify these sections, and we hope you agree.

**Highlights:**

- Voices elicit impressions about their naturalness, which affect interactions between humans as well as with artificial agents
- Despite its intuitive appeal and practical importance, a systematic understanding of voice naturalness is elusive – the concept is scientifically ill-defined
- Current voice naturalness research is situated within different research domains that resemble echo chambers within science – they neither cross-refer to one another nor to current voice perception theory
- This paper offers a concise conceptual framework by proposing a taxonomy with two distinct types: deviation-based naturalness and human-likeness-based naturalness
- This is compiled into practical recommendations and perspectives for naturalness research, because in a world of digital agents, understanding the determinants for how humans perceive naturalness in social stimuli is a priority

Click here to view linked References

# Understanding Voice Naturalness

Christine Nussbaum<sup>1,2,6</sup>, Sascha Fröhholz<sup>3,4,6</sup>, and Stefan R. Schweinberger<sup>1,2,5,6,7</sup>

<sup>1</sup>Department for General Psychology and Cognitive Neuroscience, Friedrich Schiller University Jena,  
07743 Jena, Germany

<sup>2</sup>Voice Research Unit, Friedrich Schiller University, 07743 Jena, Germany

<sup>3</sup>Department of Psychology, University of Oslo, 0371 Oslo, Norway

<sup>4</sup>Cognitive and Affective Neuroscience Unit, University of Zurich, 8050 Zurich, Switzerland

<sup>5</sup>Swiss Center for Affective Sciences, University of Geneva, 1222 Geneva, Switzerland

<sup>6</sup>The Voice Communication Sciences (VoCS) MSCA Doctoral Network

<sup>7</sup>German Center for Mental Health (DZPG), Site Jena-Halle-Magdeburg, Germany

Correspondence should be addressed to Christine Nussbaum (<https://www.allgpsy.uni-jena.de/christine-nussbaum/>), Department for General Psychology and Cognitive Neuroscience, Friedrich Schiller University Jena, Am Steiger 3/1, 07743 Jena, Germany. Tel: +49 (0) 3641 945934, E-Mail: [christine.nussbaum@uni-jena.de](mailto:christine.nussbaum@uni-jena.de). Supplemental materials to this work are accessible on the associated OSF-repository: [https://osf.io/asfqv/?view\\_only=62f8d88705bb4363903983c8bd08a2cf](https://osf.io/asfqv/?view_only=62f8d88705bb4363903983c8bd08a2cf)

1  
2  
**Abstract**  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24

Perceived naturalness of a voice is a prominent property emerging from vocal sounds, which affects our interaction with both human and artificial agents. Despite its importance, a systematic understanding of voice naturalness is elusive. This is due to (a) conceptual underspecification, (b) heterogeneous operationalization, (c) lack of exchange between research on human and synthetic voices, and (d) insufficient anchoring in voice perception theory. This review reflects on current insights into voice naturalness by pooling evidence from a wider interdisciplinary literature. Against that backdrop, it offers a concise definition of naturalness and proposes a conceptual framework rooted both in empirical findings and theoretical models. Finally, it identifies gaps in current understanding of voice naturalness and sketches perspectives for empirical progress.

**Keywords:** Naturalness, Human-likeness, Voice perception, Authenticity, Voice synthesis

25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## Naturalness: a prominent aspect of voice perception

Naturalness plays a significant role in how we perceive our environment through sight, sounds, taste, and touch. For example, perceptions of naturalness influence food choices, environmental preferences, and social trust [1–3]. From a biological perspective, perceptions of naturalness may be considered an adaptive norm, where behaviors or traits that significantly deviate from this norm are considered “unnatural”. Beyond the biological context, the recent emergence of AI-generated digital and virtual contexts has brought human-machine interactions to everyday life, thus bringing questions about naturalness to the forefront of scientific research. One of the prime channels for communicative interactions is the voice [4], both in a purely human context and beyond – with current **voice synthesis** (see Glossary) technology quickly invading everyday life, both in good use (e.g., in customer service calls, public transport, gaming, or support platforms [5,6]) and abuse (e.g., **deepfakes** [7]).

When we hear voices, we form intuitive impressions about them within just a few hundred milliseconds [8–10]. Crucially, listeners are very sensitive to impressions of voice (un)naturalness. Unnatural voices may sound nasal or robotic, or may differ from the norm in pitch contour, temporal structure, or spectral composition; in short, there are many ways in which a voice can lack naturalness [11]. Importantly, variations in naturalness affect communicative quality [12,13]. Evidence from speech-language pathologies suggests that individuals with compromised speech naturalness are often perceived as withdrawn, cold, introverted or bored [14], which can lead to social isolation and reduced quality of life [15–17] even when speech intelligibility is preserved [18]. Accordingly, voice naturalness is a key target of speech therapy, across various voice alterations [18–20]. A recent survey on personalized speech synthesis for people who lost their biological voice found that a majority prefers a more natural-sounding voice, even at the cost of some loss in intelligibility, both as users and listeners [21]. Thus, for human-to-human interaction, reduced voice naturalness consistently has negative implications.

However, this is less clear for human-machine interaction. The Computers-Are-Social-Actors (CASA) framework proposed in the 1990s [22] assumed that we treat artificial agents like humans, fueling an (implicit) naturalness-is-better bias. This spurred efforts to create synthetic voices that resemble human vocal expression [23,24], even when the link between naturalness and success in human-machine interactions remains far from fully understood. While initial findings suggested that reduced naturalness in synthetic voices compromises likeability, trustworthiness, and pleasantness [11,25–28], contemporary synthetic voice design questions a “one size fits all” idea and instead advocates solutions tailored to specific applications [29]. Accordingly, maximum human-likeness of synthetic voices may not always be required or desirable. Instead, synthetic voice preferences may depend on the features of the listeners [27,30], the device [31–33], and its specific function [6,25,31]. Understanding and incorporating such preferences seems crucial for the success and acceptance of these devices [28].

Given its widespread practical importance, the role of voice naturalness warrants scientific scrutiny. However, although many recent studies provide useful empirical insights, the current landscape resembles a patchwork rather than a cohesive research field. There are four key issues within the existing literature: (a) conceptual underspecification, (b) heterogeneous operationalization, (c) lack of exchange between research domains, and (d) insufficient anchoring in voice perception theory. These challenges have likely precluded a systematic understanding of vocal naturalness, limited visibility to a wider audience, obscured crucial research questions, and led to a divergence between theory and practice. The following sections elaborate on each of these problems, before proposing concrete measures to address them.

## Current problems in voice naturalness research

### *Conceptual underspecification*

Voice naturalness lacks a consistent definition and terminology in the literature (see **Figure 1A-B**). Many papers do not even provide an explicit definition of naturalness (see **Box 1**). In these studies, the conceptualization of naturalness must be inferred from the empirical design. If definitions are

provided, they often vary across research contexts (see **Table 1** for examples). In speech-language pathology, some researchers refer to the definition provided by Yorkston and colleagues (1999): “Naturalness is defined as conforming to the listener’s standards of rate, rhythm, intonation, and stress patterning and to the syntactic structure of the utterance being produced” [17,34]. In contrast, research on synthetic and non-human voices usually defines naturalness as “speech most closely perceived as a human voice” [35] or “the degree to which a user feels a certain technology or system is human-like” [36]. Accordingly, many studies using synthetic voices do not refer to naturalness but to human-likeness or **anthropomorphism** of voices.

Interestingly, these definitions seem to share two important assumptions: First, that voice naturalness is a perceptual and subjective measure [37]. Second, that listeners’ naturalness perception is the result of a complex multifactorial impression formation, presumably based on the integration and weighting of many **acoustic cues** [38]. Beyond this, conceptualizations are highly heterogeneous because they are tailored to the respective empirical focus. These prevailing inconsistencies alongside heterogeneous terminology (discussed next) make it challenging to compare and integrate different insights. Therefore, there is a strong need to unite them under a concise conceptual framework.

### *Heterogeneous operationalization*

A common consequence of inconsistent conceptualization is heterogeneous operationalization. Primarily, this concerns the studied vocal categories and features, which include human vs. synthetic voices [30,39–42]; cartoon voices [43]; pathological voices such as in individuals with Parkinson’s disease [44–47], **tracheoesophageal speech** [48,49], **dysarthria** [50–53], Down syndrome [54], or stuttering [19]; acoustically manipulated human voices [55]; vocal fry [56]; as well as different accents [57,58], dialects [59], age groups [60–62], and gender identities [20,63,64]. In addition, it concerns the experimental designs and measurements, especially rating scales which differ in the number of levels and denominations of endpoints. For example, in one study participants were asked “How natural is the audio?” from “1 – natural” to “5 – unnatural” [65], in another one they rated voices on a 10-point-scale from “very natural, human-like” to “very mechanical, robot-like” [58], or made a binary

classification of voices as either human or computer-generated [37]. In principle, such empirical heterogeneity can be a powerful source of insight, potentially revealing the degree to which methodological aspects affect results. For example, there is recent evidence from face perception that differences in rating scales may not have a large impact on outcome [66]. However, it cannot be concluded that this generalizes to naturalness ratings, and the insufficient report of empirical details impedes a meaningful comparison of findings. Specifically, it is often not stated how naturalness and the related experimental task were explained to the listeners – but instructions can be crucial determinants of study outcome. Further, the precise acoustic properties of voice material often remain elusive, bearing a risk for potential undetected confounds. Finally, only few studies provide measurements on interrater reliability [67]. To help address these issues, **Box 2** provides a compilation of practical recommendations as guidance for future research.

### *Lack of exchange between different research domains*

Research on voice naturalness is inherently interdisciplinary, with two main domains: speech-language pathology and synthetic voices. However, while the scientific findings are acknowledged and referenced within each domain, these domains are poorly interconnected. **Figure 1C** illustrates this via a cross-citation analysis using VOSViewer [68], showing several distinct clusters of studies reminiscent of echo chambers which are frequently discussed in social media [69]. Poor interconnectivity is not unique to naturalness but can affect many other research domains within person perception. Consider fields with different research traditions, such as impression formation according to social psychological models of inter-group perception versus face/voice perception models. These models were developed for different types of perceptual cues, and different two-factor models with different labels have been proposed in both cases (e.g., warmth vs. competence, e.g. [70]; or trustworthiness vs. dominance, e.g. [71]). More recently though, these fields arguably benefited from interconnectivity, with substantial research to link these distinct clusters and uncover both these specific taxonomies and their empirical relationships [72,73]. In the case of voice naturalness, however, two recent systematic literature reviews on pathological [17] and synthetic voices [23] do not have a single reference in common. One

1 might argue that this is not problematic, because the different disciplines simply have different  
2 interests and readerships. However, some intriguing commonalities and systematic patterns only  
3 emerge when pooling evidence from all available angles. For example, across synthetic, pathological,  
4 and acoustically manipulated voices, converging evidence emerges for a strong effect of pitch variation  
5 on perceived naturalness [14,26,74]. Further, while several studies failed to find an **uncanny valley**  
6 [75] effect for synthetic voices [11,76], a recent study suggests it might exist for pathological voices  
7 [77]. This lack of exchange between research fields has not only precluded relevant insights but has  
8 likely impeded the visibility and impact of voice naturalness research as a whole.  
9  
10  
11  
12  
13  
14  
15  
16  
17

### *18           19           20           21           22           23           24           25           26           27           28           29           30           31           32           33           34           35           36           37           38           39           40           41           42           43           44           45           46           47           48           49           50           51           52           53           54           55           56           57           58           59           60           61           62           63           64           65*

### *Insufficient anchoring in voice perception theory*

The majority of naturalness research comes from applied fields, aiming to optimize artificial agents or  
to improve the quality of life in patients with voice disorders. These findings provide valuable practical  
knowledge, but they are insufficiently anchored in voice perception theory. As an illustration, we  
added ten influential, theory-building voice perception publications to the VOSViewer analysis (**Figure 1C**), with the outcome suggesting that these tend to be ignored by most previous naturalness research.  
Indeed, several authors have pointed out that research on voice naturalness is lacking theoretical  
perspectives on voice perception and voice analysis [17,23]. This leaves us with an intriguing  
divergence between increasing applied knowledge in rapidly developing branches (especially synthetic  
voices) on the one hand, and a simultaneous lack of understanding of basic mechanisms on the other  
hand. To fully understand how naturalness affects our perception and response to voices, this void  
needs to be filled.

### *Towards a concise framework for voice naturalness*

After identifying key problems that impede a systematic understanding of naturalness in voices, a  
logical next step is to propose concrete measures to address them, starting with a conceptual  
framework for the explicit definition of naturalness in voices.

*Definitions of naturalness*

We propose a taxonomy with two distinct types: Deviation-based naturalness and human-likeness-based naturalness (**Figure 2**). In **deviation-based naturalness**, naturalness is defined as the deviation from a reference that represents maximum naturalness. Example instructions for raters could be “Does this voice sound distorted?”, “Does this voice sound unusual?”, or just “Does this voice sound natural?”. This conceptualization needs two important specifications: the reference representing maximum naturalness, and the type of deviation. In some cases, the reference is explicitly provided e.g. through a comparison or baseline stimulus (see [78]). However, in many studies, raters are instructed to use an inner implicit reference that is based on their experience and expectations, e.g., judge whether “it conforms to the expected standard of unimpaired speech” [52]. The type of deviation is specified through the vocal material. It can virtually cover all acoustic features, ranging from specific manipulations (e.g., spectral features or speech rate [79–81]) to complex multivariate vocal patterns (e.g., in distorted or pathological voices [82]).

**Human-likeness-based naturalness** defines naturalness by its resemblance to a real human voice. Instructions for raters could be “Does this voice sound like a real human speaker?” or “How human-like does the voice sound to you?” Compared to the deviation-based definition, the concept of human-likeness-based naturalness requires an additional obligatory assumption: the existence of a non-human voice space. This highlights the notion of a categorical boundary to human voices, although the transition between categories can be continuous. In other words, a definition of human-likeness is only meaningful if we assume that voices can be non-human in principle. Although deviation-based naturalness may, in certain cases, cross the boundary to the non-human voice space, this boundary is not essential for its definition. Apart from this critical distinction, however, human-likeness-based naturalness may represent a special case of deviation-based naturalness: the reference is a human voice (or listeners’ representation of a human voice), and the deviation is assessed along the human/non-human spectrum. The above considerations suggest that the human-likeness-based conceptualization is particularly well-suited for research into synthetic voices.

With this taxonomy, we provide a flexible and intuitive reference for the explicit definition of naturalness alongside its underlying assumptions. With future research committed to one conceptual framework, systematic integration and comparison of findings could be greatly facilitated. In fact, both conceptualizations seem already prevalent (see **Table 1**), but often remain implicit through certain design choices only (see **Box 1**). For example, comparing human to synthetic voices typically implies human-likeness-based naturalness, whereas assessment of pathological voices often employs the deviation-based approach. One study deserves particular mention: Diel and Lewis [77] studied the uncanny valley effect in different types of unnatural voices. They found that impressions of uncanniness resulted from “deviation from familiar categories” rather than “categorical ambiguity”. This could reflect initial empirical observations in line with our proposed conceptual distinction.

#### *Delimiting distinctiveness and authenticity*

The following section briefly discusses the demarcation of the proposed definitions of naturalness from two established concepts in perception research, starting with distinctiveness. Distinctiveness, as opposed to typicality, has been defined as the degree to which faces or voices stick out due to rare or unusual features, and this concept is commonly used to refer to identity [83,84]. According to face or voice space models, individual instances are represented along multiple perceptual dimensions, and they appear distinctive if they deviate substantially from a central tendency or norm in that space. Our deviation-based definition of naturalness is closely related to the concept of distinctiveness, as both share two critical features: a norm/reference and a deviation. However, distinctiveness, as a different concept, can capture multiple forms of deviations beyond naturalness. Accordingly, while unnatural voices would commonly be perceived as somewhat distinctive, natural voices can be distinct or typical. However, one may speculate that impressions of human-based naturalness could be quite independent from impressions of distinctiveness under certain conditions. For instance, a person who is very accustomed to a smart-speaker device may not rate synthetic voices as very distinctive but still clearly non-human. In that vein, the link between distinctiveness and naturalness may not primarily be a conceptual but an empirical matter, requiring future inspection.

A second concept that deserves particular consideration is authenticity. In the scientific literature, authenticity is an established term with meaning that may refer to vocal emotion, identity or gender – rather than the holistic impression of a voice. Emotional authenticity, for example, refers to the distinction between a posed and a “real”/spontaneous emotional expression, which leads to differential behavioral and neural outcomes [85–87]. In the context of voice cloning and the now very prevalent challenge of deepfakes [7], identity authenticity is assessed with regard to a specific speaker. In principle, authenticity can be assessed with regard to manifold social signals, including age, gender, or even personality [88,89]. At first sight, the concepts of authenticity and naturalness appear highly similar. In fact, when ChatGPT was prompted for synonyms of naturalness, authenticity was its first reply (**Figure 1B**), which may suggest that in openly accessible online sources, these two terms are indeed frequently occurring in an interchangeable manner. Accordingly, it might be argued that authenticity is just a special form of deviation-based naturalness, with a more specific reference. E.g. “Does this sound like a natural voice?” is converted into “does this sound like a natural emotional expression?”. However, if considered against the backdrop of voice perception theory, it becomes apparent that assessments of naturalness and authenticity appear at different stages of voice processing (see **Figure 3**). Thus, it would be preferable to keep the concepts of naturalness and authenticity rather separate.

## Converging evidence

In our view, understanding voice naturalness requires pooling evidence from all relevant fields. Even when these may nurture different perspectives on voice naturalness, they are united by overarching questions: How do we form an impression about voice naturalness? Which acoustic features affect this impression? How does naturalness impact perception, interaction, and communication? Can we understand differences across individuals and listening contexts?

In principle, conceptual progress for disintegrated – but also highly interdisciplinary – naturalness research can be achieved by two measures: (a) converting empirical heterogeneity from an

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
impediment into an advantage and (b) fostering mutually beneficial exchange between fields.

Awareness of the interdisciplinary nature of the field is crucial for implementing both measures: First, publications need to be findable and accessible, preferably through the establishment of common terminology that feeds into common keywords. Second, findings need to be communicated inclusively for readerships from diverse backgrounds. Finally, conceptual and empirical aspects need to be reported with sufficient detail to promote comparability. In **Box 2**, these suggestions were converted into practical recommendations.

Progress along these lines will not only enhance mutual inspiration between clinicians and engineers but could also foster innovative health technology. For instance, voice naturalness is a key objective for cochlear implant (CI) research, where a sensory prosthesis restores hearing in people with sensorineural deafness by resynthesizing auditory signals for direct electrical stimulation of the cochlea [90], and real-time synthesis in CI sound processors could be modified to achieve better perceptual outcomes, ultimately benefitting quality of life [91]. For people who are predicted to lose their personal voice due to progressive disorders such as ALS or due to planned **laryngectomy**, current voice banking technology already allows for personalized speech synthesis with the patient's former individual voice, often with remarkably high ratings of both naturalness and authenticity [21,92].

## Naturalness research rooted in voice perception theory

So far, no considerable efforts have been made to link naturalness perception to distinct stages of voice processing. As discussed earlier, the topic of voice naturalness is highly influenced by research perspectives from applied sciences and seemingly less by basic voice research and its theoretical approaches. However, neurocognitive models of voice perception can provide process-related perspectives on multi-level voice perception and voice information analysis. This allows us to link the mechanisms underlying voice naturalness assessments to the appropriate level of voice analysis. Influential theories of voice perception propose sequential and partly hierarchical stages of voice processing, including a major distinction between mechanisms for voice object analysis (i.e. perception

of an auditory stimulus as a voice) as initial stages that are followed by the analysis of communicative  
1  
2 and social content carried by the voice signal [4,93–95].  
3  
4

This processing distinction between voice object analysis and voice content analysis is relevant to the  
5 conceptual distinction between the assessment of voice naturalness on the one hand and the  
6 assessment of the authenticity of expressed voice content on the other hand (**Figure 3**). Assessing the  
7 naturalness of voices is conceptually associated with the initial levels of voice object analysis, including  
8 the stages of low-level auditory analysis and the analysis of structural voice patterns. Humans  
9 presumably assess acoustic feature deviations and acoustic feature likeness as low-level naturalness  
10 assessments [96], whereas assessing pattern deviations and pattern likeness concerns the assessments  
11 of natural or unnatural spectrotemporal voice profiles [97].  
12  
13

Whereas voice naturalness assessments likely take place at the earlier stages of voice object analysis,  
14 authenticity assessments likely take place at later stages involving voice information analysis. Voices  
15 are used as carriers to express communicative and social content. For example, voices are used for  
16 speech communication, emotional expressions, and to produce individual voice characteristics. Such  
17 voice content could be either spontaneous and authentic, or it could be acted and thus rather  
18 nonauthentic [98]. This authentic/non-authentic distinction specifically also concerns person-specific  
19 identity information in voices, which could be real or fake [7]. Such authenticity assessments might be  
20 independent of naturalness assessments, although there is also a possibility of mutual influences. For  
21 instance, perceiving a voice as unnatural might bias non-authenticity judgments of voice content, and  
22 vice versa.  
23  
24

## Perspectives for future research

Our theoretical considerations on the processing of voice naturalness call for investigations of its time-  
25 course and underlying brain mechanisms – relative to authenticity assessment but also to other voice  
26 characteristics. Initial evidence suggests that voice naturalness affects the brain response as early as  
27 200 ms after voice onset and interacts with the processing of vocal emotions [99–101]. However,  
28  
29

comparably early effects have been found for authenticity assessments [86,102,103]. Although the  
1 interpretability of these findings is limited due to the potential influence of acoustic confounds, they  
2 suggest that naturalness and authenticity assessments both are fast and fundamental parts of voice  
3 perception. However, electrophysiological insights directly comparing the time-course of naturalness  
4 and authenticity are elusive, as is their interplay with impressions of age, gender, or personality traits.  
5 A recent EEG study suggests that many first impressions formed from voices are highly intercorrelated  
6 [8], but for naturalness we are currently limited to behavioral data that point towards interactions with  
7 age, gender, and emotion perception [60,63,74]. In a broad sense, naturalness impressions are always  
8 formed against a specific context, whether that context refers to the voice itself or the properties of  
9 the interaction. Accordingly, whether the same voice is assessed in an all-human or human-machine  
10 interaction context could make a crucial difference.

In that vein, while this article focuses on understanding naturalness in voices from an interdisciplinary  
26 perspective, we wish to emphasize the multisensory perspective of naturalness research. In fact,  
27 substantial research in the domain of faces has compared the perceived naturalness or realism of  
28 synthesized versus real faces (for a systematic review and meta-analysis, see [104]). Recent research  
29 even demonstrated conditions in which synthesized faces can be perceived as more human than  
30 genuine human faces. Moreover, an attempt to identify the visual features that trigger such a  
31 paradoxical facial “hyperrealism” effect suggested contributions of typicality, familiarity,  
32 attractiveness and low memorability [105]. Although this interpretation was based on qualitative  
33 reports and requires converging evidence, such research can inspire the systematic search for  
34 commonalities or differences between mechanisms that trigger judgments of voice or face  
35 naturalness. Ultimately, naturalness research should also systematically consider interactions between  
36 vocal and visual aspects of naturalness in combination. Indeed, accumulating evidence suggests a  
37 complex interplay of visual appearance, vocal features, behavior and the interactional context for the  
38 acceptance of virtual agents [28,31–33,106–113].

Beyond humans, vocalizations are abundant in the animal kingdom. Many animals can manipulate and adapt their vocal calls to specific situations or needs. For instance, birds living in urban environments modify their song in frequency or amplitude, to avoid masking by constant anthropogenic noise [114]. While this reduces the risk of not being heard by conspecifics, the degree to which such urban-induced changes to natural patterns of vocalization may have other consequences to communication seems unclear at present. Potentially, with appropriate adaptations, the present taxonomy could be useful to promote an understanding of animal voice naturalness as well.

Finally, very recent fMRI research has uncovered a cortical-striatal brain network that is involved when listeners try to distinguish deepfake from real speaker identities [7]. Such research is relevant also because the accelerating spread of misinformation via social media is now considered a major problem which compromises societal cohesion [69,115]. While large-scale misinformation is still mostly text-based as of today, next-generation deepfakes likely will be even more efficient vehicles of misinformation. This is because they efficiently instrumentalize person-related trust via high-level perceptual deception. From that perspective, better understanding of characteristics of “successful” vocal deepfakes and their processing in the brain may be one important component for strengthening human resilience to fake information of the future.

## Concluding remarks

Naturalness in voices is a highly intuitive concept, but one that is scientifically underspecified and far from systematically understood, despite considerable research efforts. To address this, we propose a conceptual framework for voice naturalness. Our taxonomy, comprised of deviation-based naturalness and human-likeness-based naturalness, is rooted in voice perception theory, and is inspired by interdisciplinary empirical findings. The new framework offers the flexibility that is necessary to be applicable across diverse empirical designs, while at the same time promoting comparability across research domains. This conceptual groundwork is complemented with several practical recommendations to bridge previously unconnected approaches and better integrate this highly

interdisciplinary field. This provides a foundation for conjoined efforts towards more systematic future research on numerous open questions on voice naturalness (see Outstanding questions). While the focus is on voices here, we ultimately opt for a multisensory perspective on naturalness research. In a world that is increasingly dominated by digitally synthesized agents, it seems important to identify the multifaceted determinants for human perception of naturalness in social stimuli.

## Figure Legends

### *Figure 1. Terminology and interconnectivity of voice naturalness research*

**A)** Word cloud depicting synonyms and closely related concepts from 72 publications that target naturalness in voices (for details, see **Box 1**). Word size represents number of occurrences. **B)** A similar word cloud but generated by ChatGPT (<https://chatgpt.com/?oai>, 29.04.2024), when prompted to generate 10 synonyms each for pathological, synthetic/manipulated, and healthy voices, together with relative occurrence frequency. The full prompt, the generated response, and a reflection on its strengths and limitations are accessible on [OSF](#)(S6 – S8). **C)** A bibliographic network visualization using VOSviewer [68], covering publications related to voice naturalness (S1) across different domains and 10 basic voice theory papers (S2). Each colored dot represents a publication and grey links represent citations. Size of the dots indicate the number of links to other publications. Clustering (depicted by different dot colors) is performed automatically in VOSviewer. Closer inspection reveals that green refers to basic voice theory papers, red corresponds predominantly to papers on pathological voices and blue refers to synthesized/manipulated voices. A full documentation and an interactive version of the bibliographic network can be found on [OSF](#) (S9 – S12).

### *Figure 2. A conceptual framework for the definition of voice naturalness*

Assessing the naturalness of voices requires a reference frame (left panel), which is most commonly represented by the voice production system of humans. This human production system sets the reference either as individual voice samples (explicit target voice) or as prototype voice representations (implicit prototype voice), against which test voice samples (right panel) are assessed for naturalness. Two types of naturalness assessments are proposed (middle panel). The deviation-based approach assesses naturalness in terms of distance away from the reference, while the human-likeness-based approach assesses naturalness according to its similarity to the reference. Deviation in voice naturalness can occur, for example, due to clinical conditions, voice manipulations, and acoustic artifacts. Human-likeness-based naturalness defines naturalness by its resemblance to a real human voice. Human-likeness can be assessed from audio samples by judging whether or not they lie within the perimeter of an acceptable human voice border.

**Figure 3. Rooting voice naturalness in voice processing theory**

Theories of voice perception suggest a multi-level processing approach for voice samples (left panel), which involves analyzing these samples based on their features and auditory object patterns (middle panel), followed by an analysis of the information conveyed by the voice signals (right panel). Assessing the naturalness of voices appears at the level of voice features (low-level auditory analysis) and voice object analysis (voice structural analysis) and includes the assessment of acoustic deviations and acoustic likeness, as well as the assessment of pattern deviations and pattern likeness to reference voice samples. Unlike naturalness assessments, authenticity judgments mainly concern the assessment of communicative and social content carried by the voice signal at the level of voice information analysis. Such voice content can be expressed either spontaneously (authentic) or can be enacted (non-authentic), or it could be real or fake in relation to person-related identity information. Naturalness and authenticity assessments may have mutual influences.

**Table 1. Examples definitions of deviation-based and human-likeness-based voice naturalness**

Conceptualization	Definition	Reference
1 Deviation-based 2 naturalness	3 "Naturalness was defined as conforming to the 4 listener's standards of rate, rhythm, intonation, 5 and stress patterning and to the syntactic 6 structure of the utterance being produced." (p. 7 4687) 8 "Speech naturalness can be described as how 9 the speech of a person with a speech disorder 10 compares with that of typical speech or, in the 11 case of an acquired disorder, how an 12 individual's speech compares to its premorbid 13 state" (p. 1134) 14 "Speech naturalness refers to a rather broad 15 perceptual impression representing the overall 16 quality of a person's speech output in relation 17 to what is conceptualized as normal or natural" 18 (p. 1633/1634) 19 "[...] degree to which individuals sound 20 'different' from healthy peers" (p. 1265)	[44] [14] [51] [53]
21 Human-likeness- 22 based naturalness	23 "Human likeness has been used [...] to describe 24 how accurately the machine is able to imitate a 25 human." (p. 2864) 26 "Naturalness refers to whether synthetic 27 speech is perceived as uniquely human, despite 28 being computer-generated." (p. 5) 29 "Natural speech is the speech most closely 30 perceived as a human voice." (p. 10)	[26] [21] [35]
31	32	33
34 Combination of 35 both	36 "Naturalness refers to how closely the output 37 sounds like human speech." (p. 389.e1) 38 "By naturalness, we understand the voice 39 stimulus to be perceived as a plausible outcome 40 of the human speech production system" (p. 1) 41 "[...] voices which sound like they could come 42 from an actual human being (which should be 43 rated as more natural) and voices that sound 44 more fictitious, such as a cartoon character or a 45 monster (which should be rated as less natural)." (p.429)	[42] [74] [57]

46 Note. Definitions are all original quotes from the respective references. The full compilation of  
47 extracted definitions can be accessed on [OSF](#) (S5). Note that the mapping of definitions to the  
48 conceptualization of naturalness was carried out by us and not the authors of the original  
49 publications.

50  
51  
52  
53  
54  
55  
56  
57  
58 **Box 1: A field in numbers**  
59  
60  
61  
62  
63  
64  
65

For a more systematic overview of scientific insights into naturalness in voices, a focused literature search on Web of Science was conducted on 26 April 2023 using the search terms “naturalness AND voice” or “human-likeness AND voice”, which was repeated on 28 May 2024 to detect the most recent papers. This initial search resulted in 339 articles, to which the following inclusion criteria were applied:

(1) Language of publication was English. (2) Papers were published in peer-reviewed journals or as a conference contribution. (3) Voice naturalness/human-likeness was either measured or manipulated.

(4) Papers reported either a quantitative empirical analysis of human performance/perception data or a literature integration of such works. Thus, works on automatic naturalness classification and mere descriptions of toolboxes or datasets were excluded. (5) Finally, the search was focused on spoken utterances, excluding singing voices and non-linguistic vocalizations. Following these criteria, the reference lists of the identified articles were also screened for relevant publications. For a full documentation of all included papers and a reflection on potential biases in the literature search, please refer to OSF (S1, S3-S5).

In total, 72 articles were identified, covering a time range from 1984 to 2024. Thirty-eight (53%) were published in the last 5 years. Sixty-seven report behavioral empirical data, of which 48 are predominantly ratings. Two are literature reviews, and three used neurophysiological measures. Regarding voice category, 33 used synthetic, 18 human-pathological, 6 human-manipulated and 5 healthy human voices. Ten used more than one of these voice categories. In only 32 papers, an explicit definition of naturalness could be identified (see Table 1 for examples and OSF (S5) for a full list). These articles presented a large variability in wording and vocabulary. In an attempt to capture this verbal space, all articles were scanned for synonyms and closely related concepts of naturalness. The output is captured in the word cloud in **Figure 1A**. Subsequently, these were compared to the articles' keywords: 58 papers provided keywords, but only 32 had keywords related to naturalness or any of its synonyms. Finally, the conceptualization of naturalness was coded according to the taxonomy we proposed. In case no definition of naturalness was provided, the ‘implicit’ conceptualization was

inferred from the research design. With this approach, we concluded that 26 employed a deviation-based conceptualization, 35 used human-likeness, and 11 used a combination of both.

**Box 2: Practical recommendations for voice naturalness research**

Research on voice naturalness is highly interdisciplinary. To make future research accessible to a wider readership across disciplines, and allow comparability and integration of findings, awareness of this interdisciplinarity is crucial. Here is a compilation of some practical recommendations as a tentative roadmap for future research:

- Offer a concise definition of voice naturalness to both participants and readers. With the taxonomy of naturalness, this paper offers a conceptual framework that can be tailored to any empirical design, e.g. by specifying the reference and the type of deviation under study. If used consistently, this taxonomy offers a quick orientation for readers and fosters comparability across findings.
- Use consistent keywords to make relevant research findable across disciplines. We recommend “naturalness”, “human-likeness” or, in appropriate cases, “authenticity”.
- Include full reports on methodological details. Specifically, this concerns acoustic manipulations that target voice naturalness, measurements (i.e. rating scales used to assess naturalness impressions), instructions to raters, and reports on reliability. For synthetic voices, be as specific as possible on synthesis methods, toolboxes and their settings, as well as any additional processing you applied.
- Wherever possible, provide stimulus examples. This is important because readers may have a clear idea of how a male vs. female voice sounds or how an angry voice differs from a happy one, but their imagination of an (un)-natural or synthetic voice could be quite vague and differ tremendously from the actual audio material. Often, direct auditory impressions can be complementary to, and more insightful than, a list of acoustic measures and descriptions. In

some cases (i.e. when very different synthesis methods were used), differences in audio material may offer a straightforward explanation for different empirical outcomes.

- Communicate findings inclusively enough for readerships from diverse backgrounds. Provide explicit definitions (e.g. for terms like “**prosody**”, “dysarthria”, or “anthropomorphism”), avoid technical jargon including abbreviations unfamiliar to other fields (e.g. synthesis algorithms, machine learning approaches, or acoustic measures), adopt scientific standards from other fields where appropriate, and discuss findings against the wider interdisciplinary literature (i.e. linking insights into pathological voices to synthetic ones and vice versa).
- Quantify naturalness whenever it could have important implications for the ecological validity of the stimulus material, even when naturalness is not the primary focus of the study. This is especially important when using acoustic manipulations which could have unintended side effects on perceived naturalness [74,116].

### **Glossary:**

**Acoustic cues:** physical and measurable features of sounds (such as voices); these may include fundamental frequency, intensity, a range of timbre cues, or temporal characteristics. Used by listeners to inform manifold impressions about voices, such as emotion, identity, age, gender or naturalness.

**Anthropomorphism:** the attribution of human characteristics, emotions, or behaviors to non-human entities

**ChatGPT:** a chatbot developed by OpenAI, based on a large language model, that generates text based on input-prompts (GPT stands for generative pre-trained transformer)

**Deepfakes:** digitally manipulated media, such as images, videos, or voice recordings, created using deep learning techniques with the goal to convincingly display the appearance of a specific individual.

1           **Deviation-based naturalness:** Conceptualization as the deviation from a reference that represents  
2           maximum naturalness.  
3  
4

5           **Dysarthria:** impairments of speech motor subsystems due to various neurological conditions such as  
6           Parkinson's disease, amyotrophic lateral sclerosis (ALS), developmental conditions, strokes, or  
7  
8           traumatic brain injury.  
9  
10

11  
12           **Human-likeness-based naturalness:** Conceptualization of naturalness by its resemblance to a real  
13           human voice.  
14  
15

16  
17           **Laryngectomy:** surgical removal of the larynx, typically in the context of larynx cancer treatment  
18  
19

20  
21           **Prosody:** Dynamic voice intonation, as expressed in pitch, loudness, timbre, and rhythm. Sometimes  
22           also referred to as voice melody.  
23  
24

25  
26           **Tracheoesophageal speech:** a method of vocalization following total laryngectomy via a  
27           tracheoesophageal prosthesis that enables speech through esophageal vibrations.  
28  
29

30  
31           **Uncanny valley:** a sudden feeling of eeriness evoked by humanoid robots that almost approach, but  
32           do not entirely reach a human-like appearance  
33  
34

35  
36           **Voice synthesis:** creation of computer-generated voices. Common methods are articulatory  
37           synthesis, concatenative synthesis, and statistical parametric synthesis, including deep learning  
38           algorithms  
39  
40

41  
42           **Acknowledgements and Funding**  
43  
44

45  
46           We thank Simone Dahmen and Fatma Bilem for their support with the literature analysis, and the  
47           members of the Jena Voice Research Unit (<https://www.voice.uni-jena.de/>) for helpful suggestions  
48           on this project.  
49  
50

The authors gratefully acknowledge the award of funding through an EU-MSCA doctoral  
1  
2 network “Voice Communication Sciences” (action 101168998, <https://www.vocs.eu.com/>).  
3  
4

CN: I dedicate this work to our stillborn son. Thanks for changing our lives.  
5  
6  
7  
8  
9  
10

11 **Declaration of interests:**  
12  
13

14 The authors declare no competing interests.  
15  
16  
17  
18  
19  
20  
21  
22

23 **References**  
24

- 25 1. Román, S. et al. (2017) The importance of food naturalness for consumers: Results of a  
26 systematic review. *Trends in Food Science & Technology* 67, 44–57. DOI:  
27 10.1016/j.tifs.2017.06.010
- 28 2. Meier, B.P. et al. (2019) Naturally better? A review of the natural-is-better bias. *Social &*  
29 *Personality Psych* 13 (8). DOI: 10.1111/spc3.12494
- 30 3. Ode, A. et al. (2009) Indicators of perceived naturalness as drivers of landscape preference.  
31 *Journal of environmental management* 90, 375–383. DOI: 10.1016/j.jenvman.2007.10.013
- 32 4. Young, A.W. et al. (2020) Face and voice perception: Understanding commonalities and  
33 differences. *Trends Cogn Sci* 24, 398–410. DOI: 10.1016/j.tics.2020.02.001
- 34 5. Rodero, E. and Lucas, I. (2023) Synthetic versus human voices in audiobooks: The human  
35 emotional intimacy effect. *New Media & Society* 25, 1746–1764. DOI:  
36 10.1177/14614448211024142
- 37 6. Rodero, E. (2017) Effectiveness, attention, and recall of human and artificial voices in an  
38 advertising story. Prosody influence and functions of voices. *Computers in Human Behavior* 77,  
39 336–346. DOI: 10.1016/j.chb.2017.08.044
- 40 7. Roswandowitz, C. et al. (2024) Cortical-striatal brain network distinguishes deepfake from real  
41 speaker identity. *Communications biology* 7, 711. DOI: 10.1038/s42003-024-06372-6
- 42 8. Lavan, N. et al. (2024) The time course of person perception from voices in the brain. *Proc Natl*  
43 *Acad Sci U S A* 121, e2318361121. DOI: 10.1073/pnas.2318361121
- 44 9. Lavan, N. (2023) How do we describe other people from voices and faces? *Cognition* 230,  
45 105253. DOI: 10.1016/j.cognition.2022.105253
- 46 10. Jiang, Z. et al. (2024) Comparison of face-based and voice-based first impressions in a Chinese  
47 sample. *Br. J. Psychol.* 115, 20–39. DOI: 10.1111/bjop.12675
- 48 11. Kühne, K. et al. (2020) The Human Takes It All: Humanlike Synthesized Voices Are Perceived as  
49 Less Eerie and More Likable. Evidence From a Subjective Ratings Study. *Frontiers in*  
50 *NeuroRobotics* 14, 1–16. DOI: 10.3389/fnbot.2020.593732
- 51 12. Ilves, M. and Surakka, V. (2013) Subjective responses to synthesised speech with lexical  
52 emotional content: the effect of the naturalness of the synthetic voice. *Behaviour & Information*  
53 *Technology* 32, 117–131. DOI: 10.1080/0144929X.2012.702285

13. Ilves, M. et al. (2011) The Effects of Emotionally Worded Synthesized Speech on the Ratings of  
1 Emotions and Voice Quality. In , pp. 588–598, Springer, Berlin, Heidelberg  
2
14. Anand, S. and Stepp, C.E. (2015) Listener Perception of Monopitch, Naturalness, and  
3 Intelligibility for Speakers With Parkinson's Disease. *J Speech Lang Hear Res* 58, 1134–1144. DOI:  
4 10.1044/2015\_JSLHR-S-14-0243  
5
- 6 Moya-Galé, G. and Levy, E.S. (2019) Parkinson's disease-associated dysarthria: prevalence,  
7 impact and management strategies. *JPRLS Volume* 9, 9–16. DOI: 10.2147/JPRLS.S168090  
8
- 9 Damico, J.S. and Ball, M.J., eds (2019) *The SAGE Encyclopedia of Human Communication Sciences*  
10 and *Disorders*, SAGE Publications, Inc  
11
- 12 Klopfenstein, M. et al. (2020) The study of speech naturalness in communication disorders: A  
13 systematic review of the literature. *Clinical Linguistics & Phonetics* 34, 327–338. DOI:  
14 10.1080/02699206.2019.1652692  
15
- 16 Frankford, S.A. et al. (2024) Contributions of Speech Timing and Articulatory Precision to Listener  
17 Perceptions of Intelligibility and Naturalness in Parkinson's Disease. *J Speech Lang Hear Res* 67,  
18 2951–2963. DOI: 10.1044/2024\_JSLHR-23-00802  
19
- 20 Euler, H.A. et al. (2021) Speech restructuring group treatment for 6-to-9-year-old children who  
21 stutter: A therapeutic trial. *Journal of communication disorders* 89, 106073. DOI:  
22 10.1016/j.jcomdis.2020.106073  
23
- 24 Hardy, T.L.D. et al. (2020) Acoustic Predictors of Gender Attribution, Masculinity-Femininity, and  
25 Vocal Naturalness Ratings Amongst Transgender and Cisgender Speakers. *Journal of Voice* 34,  
26 300.e11-300.e26. DOI: 10.1016/j.jvoice.2018.10.002  
27
- 28 Hyppa-Martin, J. et al. (2024) A large-scale comparison of two voice synthesis techniques on  
29 intelligibility, naturalness, preferences, and attitudes toward voices banked by individuals with  
30 amyotrophic lateral sclerosis. *Augmentative and Alternative Communication* 40, 31–45. DOI:  
31 10.1080/07434618.2023.2262032  
32
- 33 Nass, C. et al. (1994) Computers are social actors. In *Proceedings of the SIGCHI conference on*  
34 *Human factors in computing systems celebrating interdependence - CHI '94*, ACM Press  
35
- 36 Seaborn, K. et al. (2021) Voice in Human-Agent Interaction. *ACM Comput. Surv.* 54, 1–43. DOI:  
37 10.1145/3386867  
38
- 39 Triantafyllopoulos, A. et al. (2023) An overview of affective speech synthesis and conversion in  
40 the deep learning era. *Proceedings of the IEEE*, 1355–1381  
41
- 42 Schreibelmayr, S. and Mara, M. (2022) Robot Voices in Daily Life: Vocal Human-Likeness and  
43 Application Context as Determinants of User Acceptance. *Frontiers in Psychology* 13, 1–17. DOI:  
44 10.3389/fpsyg.2022.787499  
45
- 46 Baird, A. et al. (2018) The Perception and Analysis of the Likeability and Human Likeness of  
47 Synthesized Speech. In *Interspeech 2018*, pp. 2863–2867, ISCA  
48
- 49 Lee, E.-J. (2010) The more humanlike, the better? How speech type and users' cognitive style  
50 affect social responses to computers. *Computers in Human Behavior* 26, 665–672. DOI:  
51 10.1016/j.chb.2010.01.003  
52
- 53 Lu, L. et al. (2021) Leveraging "human-likeness" of robotic service at restaurants. *International*  
54 *Journal of Hospitality Management* 94, 1–9. DOI: 10.1016/j.ijhm.2020.102823  
55
- 56 Cambre, J. and Kulkarni, C. (2019) One Voice Fits All? *Proc. ACM Hum.-Comput. Interact.* 3, 1–19.  
57 DOI: 10.1145/3359325  
58
- 59 Eyssel, F. et al. (2012) 'If you sound like me, you must be more human'. In *HRI' 12. Proceedings of*  
60 *the seventh annual ACM/IEEE Conference on Human-Robot Interaction : March 5-8, 2012 Boston,*  
61 *Massachusetts, USA* (Yanco, H. et al., eds), pp. 125–126, Association for Computing Machinery  
62
- 63
- 64
- 65

- 1           31. Im, H. et al. (2023) Let voice assistants sound like a machine: Voice and task type effects on  
2           perceived fluency, competence, and consumer attitude. *Computers in Human Behavior* 145,  
3           107791. DOI: 10.1016/j.chb.2023.107791
- 4           32. McGinn, C. and Torre, I. (2019 - 2019) Can you Tell the Robot by the Voice? An Exploratory Study  
5           on the Role of Voice in the Perception of Robots. In *2019 14th ACM/IEEE International  
6           Conference on Human-Robot Interaction (HRI)*, pp. 211–221, IEEE
- 7           33. Mitchell, W.J. et al. (2011) A mismatch in the human realism of face and voice produces an  
8           uncanny valley. *i-Perception* 2, 10–12. DOI: 10.1068/i0415
- 9           34. Yorkston, K.M. et al. (1999) *Management of motor speech disorders in children and adults*, Pro-  
10           ed Austin, TX
- 11           35. Mawalim, C.O. et al. (2022) Speaker anonymization by modifying fundamental frequency and x-  
12           vector singular value. *Computer Speech & Language* 73, 1–17. DOI: 10.1016/j.csl.2021.101326
- 13           36. Hu, P. et al. (2021) Dual humanness and trust in conversational AI: A person-centered approach.  
14           *Computers in Human Behavior* 119, 106727. DOI: 10.1016/j.chb.2021.106727
- 15           37. Nusbaum, H.C. et al. (1997) Measuring the naturalness of synthetic speech. *International Journal  
16           of Speech Technology* 2, 7–19
- 17           38. Mayo, C. et al. (2011) Listeners' weighting of acoustic cues to synthetic speech naturalness: A  
18           multidimensional scaling analysis. *Speech Commun* 53, 311–326. DOI:  
19           10.1016/j.specom.2010.10.003
- 20           39. Abdulrahman, A. and Richards, D. (2022) Is Natural Necessary? Human Voice versus Synthetic  
21           Voice for Intelligent Virtual Agents. *MTI* 6, 51. DOI: 10.3390/mti6070051
- 22           40. Urakami, J. et al. (2020) The Effect of Naturalness of Voice and Empathic Responses on  
23           Enjoyment, Attitudes and Motivation for Interacting with a Voice User Interface. In *Human-  
24           Computer Interaction. Multimodal and Natural Interaction* (Kurosu, M., ed), pp. 244–259,  
25           Springer International Publishing
- 26           41. Velner, E. et al. (2020) Intonation in Robot Speech. In *Proceedings of the 2020 ACM/IEEE  
27           International Conference on Human-Robot Interaction* (Belpaeme, T. et al., eds), pp. 569–578,  
28           ACM
- 29           42. Yamasaki, R. et al. (2017) Perturbation Measurements on the Degree of Naturalness of  
30           Synthesized Vowels. *Journal of Voice* 31, 389.e1-389.e8. DOI: 10.1016/j.jvoice.2016.09.020
- 31           43. Ko, S. et al. (2023) The Effects of Robot Voices and Appearances on Users' Emotion Recognition  
32           and Subjective Perception. *Int. J. Human. Robot.* 20. DOI: 10.1142/S0219843623500019
- 33           44. Abur, D. et al. (2021) Feedback and Feedforward Auditory-Motor Processes for Voice and  
34           Articulation in Parkinson's Disease. *J Speech Lang Hear Res* 64, 4682–4694. DOI:  
35           10.1044/2021\_JSLHR-21-00153
- 36           45. Klopfenstein, M. (2015) Relationship between acoustic measures and speech naturalness ratings  
37           in Parkinson's disease: A within-speaker approach. *Clinical Linguistics & Phonetics* 29, 938–954.  
38           DOI: 10.3109/02699206.2015.1081293
- 39           46. Klopfenstein, M. (2016) Speech naturalness ratings and perceptual correlates of highly natural  
40           and unnatural speech in hypokinetic dysarthria secondary to Parkinson's disease. *JRCD* 7, 123–  
41           146. DOI: 10.1558/jircd.v7i1.27932
- 42           47. Moya-Galé, G. et al. (2024) Perceptual consequences of online group speech treatment for  
43           individuals with Parkinson's disease: A pilot study case series. *International Journal of Speech-  
44           Language Pathology*, 1–16. DOI: 10.1080/17549507.2024.2330538
- 45           48. Eadie, T.L. and Doyle, P.C. (2002) Direct Magnitude Estimation and Interval Scaling of  
46           Naturalness and Severity in Tracheoesophageal (TE) Speakers. *J Speech Lang Hear Res* 45, 1088–  
47           1096. DOI: 10.1044/1092-4388(2002/087)

- 1           49. Eadie, T.L. et al. (2008) Influence of speaker gender on listener judgments of tracheoesophageal  
2           speech. *Journal of Voice* 22, 43–57. DOI: 10.1016/j.jvoice.2006.08.008  
3           50. Yorkston, K.M. et al. (1990) The effect of rate control on the intelligibility and naturalness of  
4           dysarthric speech. *The Journal of speech and hearing disorders* 55, 550–560. DOI:  
5           10.1044/jshd.5503.550  
6           51. Schölderle, T. et al. (2023) Speech Naturalness in the Assessment of Childhood Dysarthria.  
7           *American Journal of Speech-language Pathology* 32, 1633–1643. DOI: 10.1044/2023\_AJSLP-23-  
8           00023  
9           52. Lehner, K. and Ziegler, W. (2022) Clinical measures of communication limitations in dysarthria  
10          assessed through crowdsourcing: specificity, sensitivity, and retest-reliability. *Clinical Linguistics  
11          & Phonetics* 36, 988–1009. DOI: 10.1080/02699206.2021.1979658  
12          53. Vogel, A.P. et al. (2019) Speech treatment improves dysarthria in multisystemic ataxia: a rater-  
13          blinded, controlled pilot-study in ARSACS. *Journal of neurology* 266, 1260–1266. DOI:  
14          10.1007/s00415-019-09258-4  
15          54. Jones, H.N. et al. (2019) Auditory-Perceptual Speech Features in Children With Down Syndrome.  
16          *American journal on intellectual and developmental disabilities* 124, 324–338. DOI:  
17          10.1352/1944-7558-124.4.324  
18          55. Assmann, P.F. et al. (2006) Effects of frequency shifts on perceived naturalness and gender  
19          information in speech. In *INTERSPEECH*  
20          56. Venkatraman, A. and Sivasankar, M.P. (2018) Continuous Vocal Fry Simulated in Laboratory  
21          Subjects: A Preliminary Report on Voice Production and Listener Ratings. *American Journal of  
22          Speech-language Pathology* 27, 1539–1545. DOI: 10.1044/2018\_AJSLP-17-0212  
23          57. Kapolowicz, M.R. et al. (2022) Effects of Spectral Envelope and Fundamental Frequency Shifts on  
24          the Perception of Foreign-Accented Speech. *Language and speech* 65, 418–443. DOI:  
25          10.1177/00238309211029679  
26          58. Tamagawa, R. et al. (2011) The Effects of Synthesized Voice Accents on User Perceptions of  
27          Robots. *Int J of Soc Robotics* 3, 253–262. DOI: 10.1007/s12369-011-0100-4  
28          59. Mackey, L.S. et al. (1997) Effect of speech dialect on speech naturalness ratings: a systematic  
29          replication of Martin, Haroldson, and Triden (1984). *J Speech Lang Hear Res* 40, 349–360. DOI:  
30          10.1044/jslhr.4002.349  
31          60. Goy, H. et al. (2016) Effects of age on speech and voice quality ratings. *The Journal of the  
32          Acoustical Society of America* 139, 1648. DOI: 10.1121/1.4945094  
33          61. Coughlin-Woods, S. et al. (2005) Ratings of speech naturalness of children ages 8–16 years.  
34          *Percept Motor Skill* 100, 295–304. DOI: 10.2466/pms.100.2.295-304  
35          62. Baird, A. et al. (2017) Perception of Paralinguistic Traits in Synthesized Voices. In *Proceedings of  
36          the 12th International Audio Mostly Conference on Augmented and Participatory Sound and  
37          Music Experiences* (Fazekas, G. et al., eds), pp. 1–5, ACM  
38          63. Merritt, B. and Bent, T. (2020) Perceptual Evaluation of Speech Naturalness in Speakers of  
39          Varying Gender Identities. *J Speech Lang Hear Res* 63, 2054–2069. DOI: 10.1044/2020\_JSLHR-19-  
40          00337  
41          64. Baird, A. et al. (2018) The Perception of Vocal Traits in Synthesized Voices: Age, Gender, and  
42          Human Likeness. *J. Audio Eng. Soc.* 66, 277–285. DOI: 10.17743/jaes.2018.0023  
43          65. Aylett, M.P. et al. (2020) Speech Synthesis for the Generation of Artificial Personality. *IEEE Trans.  
44          Affective Comput.* 11, 361–372. DOI: 10.1109/TAFFC.2017.2763134  
45          66. Kramer, R.S.S. et al. (2024) The psychometrics of rating facial attractiveness using different  
46          response scales. *Perception* 53, 645–660. DOI: 10.1177/03010066241256221  
47          67. Martin, R.R. et al. (1984) Stuttering and speech naturalness. *The Journal of speech and hearing  
48          disorders* 49, 53–58. DOI: 10.1044/jshd.4901.53

- 1           68. van Eck, N.J. and Waltman, L. (2010) Software survey: VOSviewer, a computer program for  
2           bibliometric mapping. *Scientometrics* 84, 523–538. DOI: 10.1007/s11192-009-0146-3  
3           69. van der Linden, S. (2023) *Foolproof: Why we fall for misinformation and how to build immunity*,  
4           WW Norton & Company.  
5           70. Fiske, S.T. (2018) Stereotype Content: Warmth and Competence Endure. *Curr Dir Psychol Sci* 27,  
6           67–73. DOI: 10.1177/0963721417738825  
7           71. Todorov, A. et al. (2008) Understanding evaluation of faces on social dimensions. *Trends Cogn Sci*  
8           12, 455–460. DOI: 10.1016/j.tics.2008.10.001  
9           72. Sutherland, C.A.M. et al. (2013) Social inferences from faces: ambient images generate a three-  
10          dimensional model. *Cognition* 127, 105–118. DOI: 10.1016/j.cognition.2012.12.001  
11          73. Sutherland, C.A.M. et al. (2016) Integrating social and facial models of person perception:  
12          Converging and diverging dimensions. *Cognition* 157, 257–267. DOI:  
13          10.1016/j.cognition.2016.09.006  
14          74. Nussbaum, C. et al. (2023) Perceived naturalness of emotional voice morphs. *Cognition &*  
15          *Emotion*, 1–17. DOI: 10.1080/02699931.2023.2200920  
16          75. Mori, M. et al. (2012) The Uncanny Valley. *IEEE Robot. Automat. Mag.* 19, 98–100. DOI:  
17          10.1109/mra.2012.2192811  
18          76. Romportl, J. (2014) Speech Synthesis and Uncanny Valley. In *Text, speech and dialogue* (Horák, A.  
19          et al., eds), pp. 595–602, Springer International Publishing  
20          77. Diel, A. and Lewis, M. (2024) Deviation from typical organic voices best explains a vocal uncanny  
21          valley. *Computers in Human Behavior Reports* 14, 100430. DOI: 10.1016/j.chbr.2024.100430  
22          78. van Prooije, T. et al. (2024) Perceptual and Acoustic Analysis of Speech in Spinocerebellar ataxia  
23          Type 1. *Cerebellum*, 112–120. DOI: 10.1007/s12311-023-01513-9  
24          79. Moore, B.C.J. and Tan, C.-T. (2003) Perceived naturalness of spectrally distorted speech and  
25          music. *The Journal of the Acoustical Society of America* 114, 408–419. DOI: 10.1121/1.1577552  
26          80. Rao M V, A. et al. (2018) Effect of source filter interaction on isolated vowel-consonant-vowel  
27          perception. *The Journal of the Acoustical Society of America* 144, EL95. DOI: 10.1121/1.5049510  
28          81. Ratcliff, A. et al. (2002) Factors influencing ratings of speech naturalness in augmentative and  
29          alternative communication. *Augmentative and Alternative Communication* 18, 11–19. DOI:  
30          10.1080/aac.18.1.11.19  
31          82. Meltzner, G.S. and Hillman, R.E. (2005) Impact of Aberrant Acoustic Properties on the Perception  
32          of Sound Quality in Electrolarynx Speech. *J Speech Lang Hear Res* 48, 766–779. DOI:  
33          10.1044/1092-4388(2005/053)  
34          83. Andics, A. et al. (2010) Neural mechanisms for voice recognition. *Neuroimage* 52, 1528–1540.  
35          DOI: 10.1016/j.neuroimage.2010.05.048  
36          84. Valentine, T. et al. (2016) Face-space: A unifying concept in face recognition research. *Q J Exp*  
37          *Psychol (Hove)* 69, 1996–2019. DOI: 10.1080/17470218.2014.990392  
38          85. Lima, C.F. et al. (2021) Authentic and posed emotional vocalizations trigger distinct facial  
39          responses. *Cortex* 141, 280–292. DOI: 10.1016/j.cortex.2021.04.015  
40          86. Sarzedas, J. et al. (2024) Blindness influences emotional authenticity perception in voices:  
41          Behavioral and ERP evidence. *Cortex* 172, 254–270. DOI: 10.1016/j.cortex.2023.11.005  
42          87. Anikin, A. and Lima, C.F. (2017) Perceptual and acoustic differences between authentic and  
43          acted nonverbal emotional vocalizations. *Q J Exp Psychol (Hove)* 71, 622–641. DOI:  
44          10.1080/17470218.2016.1270976  
45          88. Kachel, S. et al. (2020) Gender (Conformity) Matters: Cross-Dimensional and Cross-Modal  
46          Associations in Sexual Orientation Perception. *Journal of Language and Social Psychology* 39, 40–  
47          66. DOI: 10.1177/0261927X19883902

- 1 89. Mills, M. et al. (2017) Expanding the evidence: Developments and innovations in clinical practice,  
2 training and competency within voice and communication therapy for trans and gender diverse  
3 people. *International Journal of Transgenderism* 18, 328–342. DOI:  
4 10.1080/15532739.2017.1329049
- 5 90. Eiff, C.I. von et al. (2022) Crossmodal benefits to vocal emotion perception in cochlear implant  
6 users. *iScience* 25, 105711. DOI: 10.1016/j.isci.2022.105711
- 7 91. Schweinberger, S.R. and Eiff, C.I. von (2022) Enhancing socio-emotional communication and  
8 quality of life in young cochlear implant recipients: Perspectives from parameter-specific  
9 morphing and caricaturing. *Frontiers in Neuroscience* 16, 956917. DOI:  
10 10.3389/fnins.2022.956917
- 11 92. Yamagishi, J. et al. (2012) Speech synthesis technologies for individuals with vocal disabilities:  
12 Voice banking and reconstruction. *Acoust. Sci. & Tech.* 33, 1–5. DOI: 10.1250/ast.33.1
- 13 93. Belin, P. et al. (2004) Thinking the voice: neural correlates of voice perception. *Trends Cogn Sci* 8,  
14 129–135. DOI: 10.1016/j.tics.2004.01.008
- 15 94. Belin, P. et al. (2011) Understanding voice perception. *Br. J. Psychol.* 102, 711–725. DOI:  
16 10.1111/j.2044-8295.2011.02041.x
- 17 95. Lavan, N. and McGettigan, C. (2023) A model for person perception from familiar and unfamiliar  
18 voices. *Commun Psychol* 1, 1–11. DOI: 10.1038/s44271-023-00001-4
- 19 96. Staib, M. and Frühholz, S. (2023) Distinct functional levels of human voice processing in the  
20 auditory cortex. *Cerebral Cortex* 33, 1170–1185. DOI: 10.1093/cercor/bhac128
- 21 97. Staib, M. and Frühholz, S. (2021) Cortical voice processing is grounded in elementary sound  
22 analyses for vocalization relevant sound patterns. *Progress in neurobiology* 200, 101982. DOI:  
23 10.1016/j.pneurobio.2020.101982
- 24 98. Pinheiro, A.P. et al. (2021) Emotional authenticity modulates affective and social trait inferences  
25 from voices. *Philosophical transactions of the Royal Society of London. Series B, Biological  
26 sciences* 376, 20200402. DOI: 10.1098/rstb.2020.0402
- 27 99. Duville, M.M. et al. (2022) Neuronal and behavioral affective perceptions of human and  
28 naturalness-reduced emotional prosodies. *Frontiers in computational neuroscience* 16, 1022787.  
29 DOI: 10.3389/fncom.2022.1022787
- 30 100. Duville, M.M. et al. (2024) Improved emotion differentiation under reduced acoustic variability  
31 of speech in autism. *BMC medicine* 22, 121. DOI: 10.1186/s12916-024-03341-y
- 32 101. Nussbaum, C. et al. (2022) Contributions of fundamental frequency and timbre to vocal emotion  
33 perception and their electrophysiological correlates. *Social Cognitive and Affective Neuroscience*  
34 17, 1145–1154. DOI: 10.1093/scan/nsac033
- 35 102. Kosilo, M. et al. (2021) The neural basis of authenticity recognition in laughter and crying.  
36 *Scientific reports* 11, 23750. DOI: 10.1038/s41598-021-03131-z
- 37 103. Conde, T. et al. (2022) The time course of emotional authenticity detection in nonverbal  
38 vocalizations. *Cortex; a journal devoted to the study of the nervous system and behavior* 151,  
39 116–132. DOI: 10.1016/j.cortex.2022.02.016
- 40 104. Miller, E.J. et al. (2023) How do people respond to computer-generated versus human faces? A  
41 systematic review and meta-analyses. *Computers in Human Behavior Reports*, 100283. DOI:  
42 10.1016/j.chbr.2023.100283
- 43 105. Miller, E.J. et al. (2023) AI Hyperrealism: Why AI Faces Are Perceived as More Real Than Human  
44 Ones. *Psychol Sci* 34, 1390–1403. DOI: 10.1177/09567976231207095
- 45 106. Cabral, J.P. et al. (2017) The Influence of Synthetic Voice on the Evaluation of a Virtual Character.  
46 In *Interspeech 2017*, pp. 229–233, ISCA
- 47 107. Ehret, J. et al. (2021) Do Prosody and Embodiment Influence the Perceived Naturalness of  
48 Conversational Agents' Speech? *ACM Trans. Appl. Percept.* 18, 1–15. DOI: 10.1145/3486580

- 108.Ferstl, Y. et al. (2021) Human or Robot? Investigating voice, appearance and gesture motion  
realism of conversational social agents. In *Proceedings of the 21th ACM International Conference  
on Intelligent Virtual Agents*, pp. 76–83, ACM
- 109.Gong, L. and Nass, C. (2007) When a Talking-Face Computer Agent is Half-Human and Half-  
Humanoid: Human Identity and Consistency Preference. *Human Comm Res* 33, 163–193. DOI:  
10.1111/j.1468-2958.2007.00295.x
- 110.Higgins, D. et al. (2022) Sympathy for the digital: Influence of synthetic voice on affinity, social  
presence and empathy for photorealistic virtual humans. *Computers & Graphics* 104, 116–128.  
DOI: 10.1016/j.cag.2022.03.009
- 111.Li, M. et al. (2023) Effects of robot gaze and voice human-likeness on users' subjective  
perception, visual attention, and cerebral activity in voice conversations. *Computers in Human  
Behavior* 141, 107645. DOI: 10.1016/j.chb.2022.107645
- 112.Parmar, D. et al. (2022) Designing Empathic Virtual Agents: Manipulating Animation, Voice,  
Rendering, and Empathy to Create Persuasive Agents. *Autonomous agents and multi-agent  
systems* 36, 1–24. DOI: 10.1007/s10458-021-09539-1
- 113.Sarigul, B. and Urgen, B.A. (2023) Audio–Visual Predictive Processing in the Perception of  
Humans and Robots. *Int J of Soc Robotics* 15, 855–865. DOI: 10.1007/s12369-023-00990-6
- 114.Lowry, H. et al. (2013) Behavioural responses of wildlife to urban environments. *Biological  
reviews of the Cambridge Philosophical Society* 88, 537–549. DOI: 10.1111/brv.12012
- 115.Kauk, J. et al. (2024) The adaptive community-response (ACR) method for collecting  
misinformation on social media. *J Big Data* 11, 1–32. DOI: 10.1186/s40537-024-00894-w
- 116.Malisz, Z. et al. (2020) Modern speech synthesis for phonetic sciences: a discussion and an  
evaluation, 487–491. DOI: 10.31234/osf.io/dxvhc

# Understanding Voice Naturalness

Christine Nussbaum<sup>1,2,6</sup>, Sascha Frühholz<sup>3,4,6</sup>, and Stefan R. Schweinberger<sup>1,2,5,6,7</sup>

<sup>1</sup>Department for General Psychology and Cognitive Neuroscience, Friedrich Schiller University Jena,  
07743 Jena, Germany

<sup>2</sup>Voice Research Unit, Friedrich Schiller University, 07743 Jena, Germany

<sup>3</sup>Department of Psychology, University of Oslo, 0371 Oslo, Norway

<sup>4</sup>Cognitive and Affective Neuroscience Unit, University of Zurich, 8050 Zurich, Switzerland

<sup>5</sup>Swiss Center for Affective Sciences, University of Geneva, 1222 Geneva, Switzerland

<sup>6</sup>The Voice Communication Sciences (VoCS) MSCA Doctoral Network

<sup>7</sup>German Center for Mental Health (DZPG), Site Jena-Halle-Magdeburg, Germany

Correspondence should be addressed to Christine Nussbaum (<https://www.allgpsy.uni-jena.de/christine-nussbaum/>), Department for General Psychology and Cognitive Neuroscience, Friedrich Schiller University Jena, Am Steiger 3/1, 07743 Jena, Germany. Tel: +49 (0) 3641 945934, E-Mail: [christine.nussbaum@uni-jena.de](mailto:christine.nussbaum@uni-jena.de). Supplemental materials to this work are accessible on the associated OSF-repository: [https://osf.io/asfqv/?view\\_only=62f8d88705bb4363903983c8bd08a2cf](https://osf.io/asfqv/?view_only=62f8d88705bb4363903983c8bd08a2cf)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Abstract**

Perceived naturalness of a voice is a prominent property emerging from vocal sounds, which affects our interaction with both human and artificial agents. Despite its importance, a systematic understanding of voice naturalness is elusive. This is due to (a) conceptual underspecification, (b) heterogeneous operationalization, (c) lack of exchange between research on human and synthetic voices, and (d) insufficient anchoring in voice perception theory. **This review** reflects on current insights into voice naturalness by pooling evidence from a wider interdisciplinary literature. Against that backdrop, **it offers** a concise definition of naturalness and proposes a conceptual framework rooted both in empirical findings and theoretical models. Finally, **it identifies** gaps in current understanding of voice naturalness and **sketches** perspectives for empirical progress.

**Keywords:** Naturalness, Human-likeness, Voice perception, Authenticity, Voice synthesis

## Naturalness: a prominent aspect of voice perception

Naturalness plays a significant role in how we perceive our environment through sight, sounds, taste, and touch. For example, perceptions of naturalness influence food choices, environmental preferences, and social trust [1–3]. From a biological perspective, perceptions of naturalness **may be considered an adaptive norm**, where behaviors or traits that significantly deviate from this norm are considered “unnatural”. Beyond the biological context, the recent emergence of AI-generated digital and virtual contexts has brought human-machine interactions to everyday life, thus bringing questions about naturalness to the forefront of scientific research. One of the prime channels for communicative interactions is the voice [4], both in a purely human context and beyond – with current **voice synthesis** (see Glossary) technology quickly invading everyday life, both in good use (e.g., in customer service calls, public transport, gaming, or support platforms [5,6]) and abuse (e.g., **deepfakes** [7]).

When we hear voices, we form intuitive impressions about them within just a few hundred milliseconds [8–10]. Crucially, listeners are very sensitive to impressions of voice (un)naturalness. Unnatural voices may sound nasal or robotic, or may differ from the norm in pitch contour, temporal structure, or spectral composition; in short, there are many ways in which a voice can lack naturalness [11]. Importantly, variations in naturalness affect communicative quality [12,13]. Evidence from speech-language pathologies suggests that individuals with compromised speech naturalness are often perceived as withdrawn, cold, introverted or bored [14], which can lead to social isolation and reduced quality of life [15–17] even when speech intelligibility is preserved [18]. Accordingly, voice naturalness is a key target of speech therapy, across various voice alterations [18–20]. A recent survey on personalized speech synthesis for people who lost their biological voice found that a majority prefers a more natural-sounding voice, even at the cost of some loss in intelligibility, both as users and listeners [21]. Thus, for human-to-human interaction, reduced voice naturalness consistently has negative implications.

1 However, this is less clear for human-machine interaction. The Computers-Are-Social-Actors (CASA)  
2 framework proposed in the 1990s [22] assumed that we treat artificial agents like humans, fueling an  
3 (implicit) naturalness-is-better bias. This spurred efforts to create synthetic voices that resemble  
4 human vocal expression [23,24], even when the link between naturalness and success in human-  
5 machine interactions remains far from fully understood. While initial findings suggested that reduced  
6 naturalness in synthetic voices compromises likeability, trustworthiness, and pleasantness [11,25–28],  
7 contemporary synthetic voice design questions a “one size fits all” idea and instead advocates solutions  
8 tailored to specific applications [29]. Accordingly, maximum human-likeness of synthetic voices may  
9 not always be required or desirable. Instead, synthetic voice preferences may depend on the features  
10 of the listeners [27,30], the device [31–33], and its specific function [6,25,31]. Understanding and  
11 incorporating such preferences seems crucial for the success and acceptance of these devices [28].  
12  
13 Given its widespread practical importance, the role of voice naturalness warrants scientific scrutiny.  
14 However, although many recent studies provide useful empirical insights, the current landscape  
15 resembles a patchwork rather than a cohesive research field. There are four key issues within the  
16 existing literature: (a) conceptual underspecification, (b) heterogeneous operationalization, (c) lack of  
17 exchange between research domains, and (d) insufficient anchoring in voice perception theory. These  
18 challenges have likely precluded a systematic understanding of vocal naturalness, limited visibility to a  
19 wider audience, obscured crucial research questions, and led to a divergence between theory and  
20 practice. The following sections elaborate on each of these problems, before proposing concrete  
21 measures to address them.

## 49 Current problems in voice naturalness research

### 52 *Conceptual underspecification*

55 Voice naturalness lacks a consistent definition and terminology in the literature (see **Figure 1A-B**).  
56 Many papers do not even provide an explicit definition of naturalness (see **Box 1**). In these studies, the  
57 conceptualization of naturalness must be inferred from the empirical design. If definitions are  
58  
59  
60  
61  
62  
63  
64  
65

provided, they often vary across research contexts (see **Table 1** for examples). In speech-language pathology, some researchers refer to the definition provided by Yorkston and colleagues (1999): “Naturalness is defined as conforming to the listener’s standards of rate, rhythm, intonation, and stress patterning and to the syntactic structure of the utterance being produced” [17,34]. In contrast, research on synthetic and non-human voices usually defines naturalness as “speech most closely perceived as a human voice” [35] or “the degree to which a user feels a certain technology or system is human-like” [36]. Accordingly, many studies using synthetic voices do not refer to naturalness but to human-likeness or **anthropomorphism** of voices.

Interestingly, these definitions seem to share two important assumptions: First, that voice naturalness is a perceptual and subjective measure [37]. Second, that listeners’ naturalness perception is the result of a complex multifactorial impression formation, presumably based on the integration and weighting of many **acoustic cues** [38]. Beyond this, conceptualizations are highly heterogeneous because they are tailored to the respective empirical focus. These prevailing inconsistencies alongside heterogeneous terminology (discussed next) make it challenging to compare and integrate different insights. **Therefore, there is a strong need to unite them under a concise conceptual framework.**

### *Heterogeneous operationalization*

A common consequence of inconsistent conceptualization is heterogeneous operationalization. Primarily, this concerns the studied vocal categories and features, which include human vs. synthetic voices [30,39–42]; cartoon voices [43]; pathological voices such as in individuals with Parkinson’s disease [44–47], **tracheoesophageal speech** [48,49], **dysarthria** [50–53], Down syndrome [54], or stuttering [19]; acoustically manipulated human voices [55]; vocal fry [56]; as well as different accents [57,58], dialects [59], age groups [60–62], and gender identities [20,63,64]. In addition, it concerns the experimental designs and measurements, especially rating scales which differ in the number of levels and denominations of endpoints. For example, in one study participants were asked “How natural is the audio?” from “1 – natural” to “5 – unnatural” [65], in another one they rated voices on a 10-point-scale from “very natural, human-like” to “very mechanical, robot-like” [58], or made a binary

classification of voices as either human or computer-generated [37]. In principle, such empirical heterogeneity can be a powerful source of insight, potentially revealing the degree to which methodological aspects affect results. For example, there is recent evidence from face perception that differences in rating scales may not have a large impact on outcome [66]. However, it cannot be concluded that this generalizes to naturalness ratings, and the insufficient report of empirical details impedes a meaningful comparison of findings. Specifically, it is often not stated how naturalness and the related experimental task were explained to the listeners – but instructions can be crucial determinants of study outcome. Further, the precise acoustic properties of voice material often remain elusive, bearing a risk for potential undetected confounds. Finally, only few studies provide measurements on interrater reliability [67]. To help address these issues, Box 2 provides a compilation of practical recommendations as guidance for future research.

### *Lack of exchange between different research domains*

Research on voice naturalness is inherently interdisciplinary, with two main domains: speech-language pathology and synthetic voices. However, while the scientific findings are acknowledged and referenced within each domain, these domains are poorly interconnected. Figure 1C illustrates this via a cross-citation analysis using VOSViewer [68], showing several distinct clusters of studies reminiscent of echo chambers which are frequently discussed in social media [69]. Poor interconnectivity is not unique to naturalness but can affect many other research domains within person perception. Consider fields with different research traditions, such as impression formation according to social psychological models of inter-group perception versus face/voice perception models. These models were developed for different types of perceptual cues, and different two-factor models with different labels have been proposed in both cases (e.g., warmth vs. competence, e.g. [70]; or trustworthiness vs. dominance, e.g. [71]). More recently though, these fields arguably benefited from interconnectivity, with substantial research to link these distinct clusters and uncover both these specific taxonomies and their empirical relationships [72,73]. In the case of voice naturalness, however, two recent systematic literature reviews on pathological [17] and synthetic voices [23] do not have a single reference in common. One

1 might argue that this is not problematic, because the different disciplines simply have different  
2 interests and readerships. However, some intriguing commonalities and systematic patterns only  
3 emerge when pooling evidence from all available angles. For example, across synthetic, pathological,  
4 and acoustically manipulated voices, converging evidence emerges for a strong effect of pitch variation  
5 on perceived naturalness [14,26,74]. Further, while several studies failed to find an **uncanny valley**  
6 [75] effect for synthetic voices [11,76], a recent study suggests it might exist for pathological voices  
7 [77]. This lack of exchange between research fields has not only precluded relevant insights but has  
8 likely impeded the visibility and impact of voice naturalness research as a whole.  
9  
10  
11  
12  
13  
14  
15  
16  
17

### *Insufficient anchoring in voice perception theory*

22 The majority of naturalness research comes from applied fields, aiming to optimize artificial agents or  
23 to improve the quality of life in patients with voice disorders. These findings provide valuable practical  
24 knowledge, but they are insufficiently anchored in voice perception theory. As an illustration, we  
25 added ten influential, theory-building voice perception publications to the VOSViewer analysis (**Figure**  
26 **1C**), with the outcome suggesting that these tend to be ignored by most previous naturalness research.  
27 Indeed, several authors have pointed out that research on voice naturalness is lacking theoretical  
28 perspectives on voice perception and voice analysis [17,23]. This leaves us with an intriguing  
29 divergence between increasing applied knowledge in rapidly developing branches (especially synthetic  
30 voices) on the one hand, and a simultaneous lack of understanding of basic mechanisms on the other  
31 hand. To fully understand how naturalness affects our perception and response to voices, this void  
32 needs to be filled.  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48

### Towards a concise framework for voice naturalness

51 After identifying key problems that impede a systematic understanding of naturalness in voices, a  
52 logical next step is to propose concrete measures to address them, starting with a conceptual  
53 framework for the explicit definition of naturalness in voices.  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

*Definitions of naturalness*

We propose a taxonomy with two distinct types: Deviation-based naturalness and human-likeness-based naturalness (**Figure 2**). In **deviation-based naturalness**, naturalness is defined as the deviation from a reference that represents maximum naturalness. Example instructions for raters could be “Does this voice sound distorted?”, “Does this voice sound unusual?”, or just “Does this voice sound natural?”. This conceptualization needs two important specifications: the reference representing maximum naturalness, and the type of deviation. In some cases, the reference is explicitly provided e.g. through a comparison or baseline stimulus (see [78]). However, in many studies, raters are instructed to use an inner implicit reference that is based on their experience and expectations, e.g., judge whether “it conforms to the expected standard of unimpaired speech” [52]. The type of deviation is specified through the vocal material. It can virtually cover all acoustic features, ranging from specific manipulations (e.g., spectral features or speech rate [79–81]) to complex multivariate vocal patterns (e.g., in distorted or pathological voices [82]).

**Human-likeness-based naturalness** defines naturalness by its resemblance to a real human voice. Instructions for raters could be “Does this voice sound like a real human speaker?” or “How human-like does the voice sound to you?” Compared to the deviation-based definition, **the concept of human-likeness-based naturalness** requires an **additional obligatory assumption**: the existence of a **non-human voice space**. **This highlights the notion of** a categorical boundary to human voices, although the transition between categories can be continuous. In other words, a definition of human-likeness is only meaningful if we assume that voices can be non-human in principle. **Although deviation-based naturalness may, in certain cases, cross the boundary to the non-human voice space, this boundary is not essential for its definition.** Apart from this **critical** distinction, however, human-likeness-based naturalness may **represent** a special case of deviation-based naturalness: the reference is a human voice (or listeners’ representation of a human voice), and the deviation is **assessed along the human/non-human spectrum**. **The above considerations suggest that the human-likeness-based conceptualization is particularly well-suited for research into synthetic voices.**

With this taxonomy, we provide a flexible and intuitive reference for the explicit definition of naturalness alongside its underlying assumptions. With future research committed to one conceptual framework, systematic integration and comparison of findings could be greatly facilitated. In fact, both conceptualizations seem already prevalent (see **Table 1**), but often remain implicit through certain design choices only (see **Box 1**). For example, comparing human to synthetic voices typically implies human-likeness-based naturalness, whereas assessment of pathological voices often employs the deviation-based approach. One study deserves particular mention: Diel and Lewis [77] studied the uncanny valley effect in different types of unnatural voices. They found that impressions of uncanniness resulted from “deviation from familiar categories” rather than “categorical ambiguity”. This could reflect initial empirical observations in line with our proposed conceptual distinction.

#### *Delimiting distinctiveness and authenticity*

The following section briefly discusses the demarcation of the proposed definitions of naturalness from two established concepts in perception research, starting with distinctiveness. Distinctiveness, as opposed to typicality, has been defined as the degree to which faces or voices stick out due to rare or unusual features, and this concept is commonly used to refer to identity [83,84]. According to face or voice space models, individual instances are represented along multiple perceptual dimensions, and they appear distinctive if they deviate substantially from a central tendency or norm in that space. Our deviation-based definition of naturalness is closely related to the concept of distinctiveness, as both share two critical features: a norm/reference and a deviation. However, distinctiveness, as a different concept, can capture multiple forms of deviations beyond naturalness. Accordingly, while unnatural voices would commonly be perceived as somewhat distinctive, natural voices can be distinct or typical. However, one may speculate that impressions of human-based naturalness could be quite independent from impressions of distinctiveness under certain conditions. For instance, a person who is very accustomed to a smart-speaker device may not rate synthetic voices as very distinctive but still clearly non-human. In that vein, the link between distinctiveness and naturalness may not primarily be a conceptual but an empirical matter, requiring future inspection.

A second concept that deserves particular consideration is authenticity. In the scientific literature, authenticity is an established term with meaning that may refer to vocal emotion, identity or gender – rather than the holistic impression of a voice. Emotional authenticity, for example, refers to the distinction between a posed and a “real”/spontaneous emotional expression, which leads to differential behavioral and neural outcomes [85–87]. In the context of voice cloning and the now very prevalent challenge of deepfakes [7], identity authenticity is assessed with regard to a specific speaker. In principle, authenticity can be assessed with regard to manifold social signals, including age, gender, or even personality [88,89]. At first sight, the concepts of authenticity and naturalness appear highly similar. In fact, when ChatGPT was prompted for synonyms of naturalness, authenticity was its first reply (Figure 1B), which may suggest that in openly accessible online sources, these two terms are indeed frequently occurring in an interchangeable manner. Accordingly, it might be argued that authenticity is just a special form of deviation-based naturalness, with a more specific reference. E.g. “Does this sound like a natural voice?” is converted into “does this sound like a natural emotional expression?”. However, if considered against the backdrop of voice perception theory, it becomes apparent that assessments of naturalness and authenticity appear at different stages of voice processing (see Figure 3). Thus, it would be preferable to keep the concepts of naturalness and authenticity rather separate.

## Converging evidence

In our view, understanding voice naturalness requires pooling evidence from all relevant fields. Even when these may nurture different perspectives on voice naturalness, they are united by overarching questions: How do we form an impression about voice naturalness? Which acoustic features affect this impression? How does naturalness impact perception, interaction, and communication? Can we understand differences across individuals and listening contexts?

In principle, conceptual progress for disintegrated – but also highly interdisciplinary – naturalness research can be achieved by two measures: (a) converting empirical heterogeneity from an

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
impediment into an advantage and (b) fostering mutually beneficial exchange between fields.

Awareness of the interdisciplinary nature of the field is crucial for implementing both measures: First, publications need to be findable and accessible, preferably through the establishment of common terminology that feeds into common keywords. Second, findings need to be communicated inclusively for readerships from diverse backgrounds. Finally, conceptual and empirical aspects need to be reported with sufficient detail to promote comparability. In **Box 2**, these suggestions were converted into practical recommendations.

Progress along these lines will not only enhance mutual inspiration between clinicians and engineers but could also foster innovative health technology. For instance, voice naturalness is a key objective for cochlear implant (CI) research, where a sensory prosthesis restores hearing in people with sensorineural deafness by resynthesizing auditory signals for direct electrical stimulation of the cochlea [90], and real-time synthesis in CI sound processors could be modified to achieve better perceptual outcomes, ultimately benefitting quality of life [91]. For people who are predicted to lose their personal voice due to progressive disorders such as ALS or due to planned **laryngectomy**, current voice banking technology already allows for personalized speech synthesis with the patient's former individual voice, often with remarkably high ratings of both naturalness and authenticity [21,92].

## Naturalness research rooted in voice perception theory

So far, no considerable efforts have been made to link naturalness perception to distinct stages of voice processing. As discussed earlier, the topic of voice naturalness is highly influenced by research perspectives from applied sciences and seemingly less by basic voice research and its theoretical approaches. However, neurocognitive models of voice perception can provide process-related perspectives on multi-level voice perception and voice information analysis. This allows us to link the mechanisms underlying voice naturalness assessments to the appropriate level of voice analysis. Influential theories of voice perception propose sequential and partly hierarchical stages of voice processing, including a major distinction between mechanisms for voice object analysis (i.e. perception

of an auditory stimulus as a voice) as initial stages that are followed by the analysis of communicative  
1  
2 and social content carried by the voice signal [4,93–95].  
3  
4

This processing distinction between voice object analysis and voice content analysis is relevant to the  
5 conceptual distinction between the assessment of voice naturalness on the one hand and the  
6 assessment of the authenticity of expressed voice content on the other hand (**Figure 3**). Assessing the  
7 naturalness of voices is conceptually associated with the initial levels of voice object analysis, including  
8 the stages of low-level auditory analysis and the analysis of structural voice patterns. Humans  
9 presumably assess acoustic feature deviations and acoustic feature likeness as low-level naturalness  
10 assessments [96], whereas assessing pattern deviations and pattern likeness concerns the assessments  
11 of natural or unnatural spectrotemporal voice profiles [97].  
12  
13

Whereas voice naturalness assessments likely take place at the earlier stages of voice object analysis,  
25  
26 authenticity assessments likely take place at later stages involving voice information analysis. Voices  
27 are used as carriers to express communicative and social content. For example, voices are used for  
28 speech communication, emotional expressions, and to produce individual voice characteristics. Such  
29 voice content could be either spontaneous and authentic, or it could be acted and thus rather  
30 nonauthentic [98]. This authentic/non-authentic distinction specifically also concerns person-specific  
31 identity information in voices, which could be real or fake [7]. Such authenticity assessments might be  
32 independent of naturalness assessments, although there is also a possibility of mutual influences. For  
33 instance, perceiving a voice as unnatural might bias non-authenticity judgments of voice content, and  
34 vice versa.  
35  
36

## Perspectives for future research

Our theoretical considerations on the processing of voice naturalness call for investigations of its time-  
53 course and underlying brain mechanisms – relative to authenticity assessment but also to other voice  
54 characteristics. Initial evidence suggests that voice naturalness affects the brain response as early as  
55 200 ms after voice onset and interacts with the processing of vocal emotions [99–101]. However,  
56  
57

comparably early effects have been found for authenticity assessments [86,102,103]. Although the  
1 interpretability of these findings is limited due to the potential influence of acoustic confounds, they  
2 suggest that naturalness and authenticity assessments both are fast and fundamental parts of voice  
3 perception. However, electrophysiological insights directly comparing the time-course of naturalness  
4 and authenticity are elusive, as is their interplay with impressions of age, gender, or personality traits.  
5 A recent EEG study suggests that many first impressions formed from voices are highly intercorrelated  
6 [8], but for naturalness we are currently limited to behavioral data that point towards interactions with  
7 age, gender, and emotion perception [60,63,74]. In a broad sense, naturalness impressions are always  
8 formed against a specific context, whether that context refers to the voice itself or the properties of  
9 the interaction. Accordingly, whether the same voice is assessed in an all-human or human-machine  
10 interaction context could make a crucial difference.

In that vein, while this article focuses on understanding naturalness in voices from an interdisciplinary  
26 perspective, we wish to emphasize the multisensory perspective of naturalness research. In fact,  
27 substantial research in the domain of faces has compared the perceived naturalness or realism of  
28 synthesized versus real faces (for a systematic review and meta-analysis, see [104]). Recent research  
29 even demonstrated conditions in which synthesized faces can be perceived as more human than  
30 genuine human faces. Moreover, an attempt to identify the visual features that trigger such a  
31 paradoxical facial “hyperrealism” effect suggested contributions of typicality, familiarity,  
32 attractiveness and low memorability [105]. Although this interpretation was based on qualitative  
33 reports and requires converging evidence, such research can inspire the systematic search for  
34 commonalities or differences between mechanisms that trigger judgments of voice or face  
35 naturalness. Ultimately, naturalness research should also systematically consider interactions between  
36 vocal and visual aspects of naturalness in combination. Indeed, accumulating evidence suggests a  
37 complex interplay of visual appearance, vocal features, behavior and the interactional context for the  
38 acceptance of virtual agents [28,31–33,106–113].

Beyond humans, vocalizations are abundant in the animal kingdom. Many animals can manipulate and adapt their vocal calls to specific situations or needs. For instance, birds living in urban environments modify their song in frequency or amplitude, to avoid masking by constant anthropogenic noise [114]. While this reduces the risk of not being heard by conspecifics, the degree to which such urban-induced changes to natural patterns of vocalization may have other consequences to communication seems unclear at present. Potentially, with appropriate adaptations, the present taxonomy could be useful to promote an understanding of animal voice naturalness as well.

Finally, very recent fMRI research has uncovered a cortical-striatal brain network that is involved when listeners try to distinguish deepfake from real speaker identities [7]. Such research is relevant also because the accelerating spread of misinformation via social media is now considered a major problem which compromises societal cohesion [69,115]. While large-scale misinformation is still mostly text-based as of today, next-generation deepfakes likely will be even more efficient vehicles of misinformation. This is because they efficiently instrumentalize person-related trust via high-level perceptual deception. From that perspective, better understanding of characteristics of “successful” vocal deepfakes and their processing in the brain may be one important component for strengthening human resilience to fake information of the future.

## Concluding remarks

Naturalness in voices is a highly intuitive concept, but one that is scientifically underspecified and far from systematically understood, despite considerable research efforts. To address this, we propose a conceptual framework for voice naturalness. Our taxonomy, comprised of deviation-based naturalness and human-likeness-based naturalness, is rooted in voice perception theory, and is inspired by interdisciplinary empirical findings. The new framework offers the flexibility that is necessary to be applicable across diverse empirical designs, while at the same time promoting comparability across research domains. This conceptual groundwork is complemented with several practical recommendations to bridge previously unconnected approaches and better integrate this highly

interdisciplinary field. This provides a foundation for conjoined efforts towards more systematic future research on numerous open questions on voice naturalness (see Outstanding questions). While the focus is on voices here, we ultimately opt for a multisensory perspective on naturalness research. In a world that is increasingly dominated by digitally synthesized agents, it seems important to identify the multifaceted determinants for human perception of naturalness in social stimuli.

## Figure Legends

### *Figure 1. Terminology and interconnectivity of voice naturalness research*

**A)** Word cloud depicting synonyms and closely related concepts from 72 publications that target naturalness in voices (for details, see **Box 1**). Word size represents number of occurrences. **B)** A similar word cloud but generated by ChatGPT (<https://chatgpt.com/?oai>, 29.04.2024), when prompted to generate 10 synonyms each for pathological, synthetic/manipulated, and healthy voices, together with relative occurrence frequency. The full prompt, the generated response, and a reflection on its strengths and limitations are accessible on [OSF\(S6 – S8\)](#). **C)** A bibliographic network visualization using VOSviewer [68], covering publications related to voice naturalness (**S1**) across different domains and 10 basic voice theory papers (**S2**). Each colored dot represents a publication and grey links represent citations. Size of the dots indicate the number of links to other publications. Clustering (depicted by different dot colors) is performed automatically in VOSviewer. Closer inspection reveals that green refers to basic voice theory papers, red corresponds predominantly to papers on pathological voices and blue refers to synthesized/manipulated voices. A full documentation and an interactive version of the bibliographic network can be found on [OSF \(S9 – S12\)](#).

### *Figure 2. A conceptual framework for the definition of voice naturalness*

Assessing the naturalness of voices requires a reference frame (left panel), which is most commonly represented by the voice production system of humans. This human production system sets the reference either as individual voice samples (explicit target voice) or as prototype voice representations (implicit prototype voice), against which test voice samples (right panel) are assessed for naturalness. Two types of naturalness assessments are proposed (middle panel). The deviation-based approach assesses naturalness in terms of distance away from the reference, while the human-likeness-based approach assesses naturalness according to its similarity to the reference. Deviation in voice naturalness can occur, for example, due to clinical conditions, voice manipulations, and acoustic artifacts. Human-likeness-based naturalness defines naturalness by its resemblance to a real human voice. **Human-likeness can be assessed from audio samples by judging whether or not they lie within the perimeter of an acceptable human voice border.**

**Figure 3. Rooting voice naturalness in voice processing theory**

Theories of voice perception suggest a multi-level processing approach for voice samples (left panel), which involves analyzing these samples based on their features and auditory object patterns (middle panel), followed by an analysis of the information conveyed by the voice signals (right panel). Assessing the naturalness of voices appears at the level of voice features (low-level auditory analysis) and voice object analysis (voice structural analysis) and includes the assessment of acoustic deviations and acoustic likeness, as well as the assessment of pattern deviations and pattern likeness to reference voice samples. Unlike naturalness assessments, authenticity judgments mainly concern the assessment of communicative and social content carried by the voice signal at the level of voice information analysis. Such voice content can be expressed either spontaneously (authentic) or can be enacted (non-authentic), **or it could be real or fake in relation to person-related identity information.**

Naturalness and authenticity assessments may have mutual influences.

**Table 1. Examples definitions of deviation-based and human-likeness-based voice naturalness**

Conceptualization	Definition	Reference
1 Deviation-based 2 naturalness	3 "Naturalness was defined as conforming to the 4 listener's standards of rate, rhythm, intonation, 5 and stress patterning and to the syntactic 6 structure of the utterance being produced." (p. 7 4687) 8 "Speech naturalness can be described as how 9 the speech of a person with a speech disorder 10 compares with that of typical speech or, in the 11 case of an acquired disorder, how an 12 individual's speech compares to its premorbid 13 state" (p. 1134) 14 "Speech naturalness refers to a rather broad 15 perceptual impression representing the overall 16 quality of a person's speech output in relation 17 to what is conceptualized as normal or natural" 18 (p. 1633/1634) 19 "[...] degree to which individuals sound 20 'different' from healthy peers" (p. 1265)	[44] [14] [51] [53]
21 Human-likeness- 22 based naturalness	23 "Human likeness has been used [...] to describe 24 how accurately the machine is able to imitate a 25 human." (p. 2864) 26 "Naturalness refers to whether synthetic 27 speech is perceived as uniquely human, despite 28 being computer-generated." (p. 5) 29 "Natural speech is the speech most closely 30 perceived as a human voice." (p. 10)	[26] [21] [35]
31	32	33
34 Combination of 35 both	36 "Naturalness refers to how closely the output 37 sounds like human speech." (p. 389.e1) 38 "By naturalness, we understand the voice 39 stimulus to be perceived as a plausible outcome 40 of the human speech production system" (p. 1) 41 "[...] voices which sound like they could come 42 from an actual human being (which should be 43 rated as more natural) and voices that sound 44 more fictitious, such as a cartoon character or a 45 monster (which should be rated as less natural)." (p.429)	[42] [74] [57]

46 Note. Definitions are all original quotes from the respective references. The full compilation of  
 47 extracted definitions can be accessed on OSF ([S5](#)). Note that the mapping of definitions to the  
 48 conceptualization of naturalness was carried out by us and not the authors of the original  
 49 publications.

50  
51  
52  
53  
54  
55  
56  
57  
58 **Box 1: A field in numbers**  
59  
60  
61

For a more systematic overview of scientific insights into naturalness in voices, a focused literature search on Web of Science **was conducted** on 26 April 2023 using the search terms “naturalness AND voice” or “human-likeness AND voice”, which was repeated on 28 May 2024 to detect the most recent papers. This initial search resulted in 339 articles, to which the following inclusion criteria **were applied**:  
(1) Language of publication was English. (2) Papers were published in peer-reviewed journals or as a conference contribution. (3) Voice naturalness/human-likeness was either measured or manipulated.  
(4) Papers reported either a quantitative empirical analysis of human performance/perception data or a literature integration of such works. Thus, works on automatic naturalness classification and mere descriptions of toolboxes or datasets **were excluded**. (5) Finally, **the search was focused** on spoken utterances, excluding singing voices and non-linguistic vocalizations. Following these criteria, the reference lists of the identified articles **were also screened** for relevant publications. For a full documentation of all included papers and a reflection on potential biases in the literature search, please refer to OSF (**S1, S3-S5**).

In total, 72 articles **were identified**, covering a time range from 1984 to 2024. Thirty-eight (53%) were published in the last 5 years. Sixty-seven report behavioral empirical data, of which 48 are predominantly ratings. Two are literature reviews, and three used neurophysiological measures. Regarding voice category, 33 used synthetic, 18 human-pathological, 6 human-manipulated and 5 healthy human voices. Ten used more than one of these voice categories. In only 32 papers, an explicit definition of naturalness **could be identified** (see Table 1 for examples and OSF (**S5**) for a full list). These articles presented a large variability in wording and vocabulary. In an attempt to capture this verbal space, all articles **were scanned** for synonyms and closely related concepts of naturalness. The output is captured in the word cloud in **Figure 1A**. Subsequently, these **were compared** to the articles' keywords: 58 papers provided keywords, but only 32 had keywords related to naturalness or any of its synonyms. Finally, the conceptualization of naturalness **was coded** according to the taxonomy we proposed. In case no definition of naturalness was provided, the ‘implicit’ conceptualization **was**

1            inferred from the research design. With this approach, we concluded that 26 employed a deviation-  
2            based conceptualization, 35 used human-likeness, and 11 used a combination of both.  
3  
4  
5  
6  
7

8            **Box 2: Practical recommendations for voice naturalness research**  
9

10            Research on voice naturalness is highly interdisciplinary. To make future research accessible to a wider  
11            readership across disciplines, and allow comparability and integration of findings, awareness of this  
12            interdisciplinarity is crucial. [Here is a compilation of](#) some practical recommendations as a tentative  
13            roadmap for future research:  
14  
15

- 16            • Offer a concise definition of voice naturalness to both participants and readers. With the  
17            taxonomy of naturalness, [this paper offers](#) a conceptual framework that can be tailored to any  
18            empirical design, e.g. by specifying the reference and the type of deviation under study. If used  
19            consistently, this taxonomy offers a quick orientation for readers and fosters comparability  
20            across findings.  
21  
22            • Use consistent keywords to make relevant research findable across disciplines. We  
23            recommend “naturalness”, “human-likeness” or, in appropriate cases, “authenticity”.  
24  
25            • Include full reports on methodological details. Specifically, this concerns acoustic  
26            manipulations that target voice naturalness, measurements (i.e. rating scales used to assess  
27            naturalness impressions), instructions to raters, and reports on reliability. For synthetic voices,  
28            be as specific as possible on synthesis methods, toolboxes and their settings, as well as any  
29            additional processing you applied.  
30  
31            • Wherever possible, provide stimulus examples. This is important because readers may have a  
32            clear idea of how a male vs. female voice sounds or how an angry voice differs from a happy  
33            one, but their imagination of an (un)-natural or synthetic voice could be quite vague and differ  
34            tremendously from the actual audio material. Often, direct auditory impressions can be  
35            complementary to, and more insightful than, a list of acoustic measures and descriptions. In  
36  
37

some cases (i.e. when very different synthesis methods were used), differences in audio material may offer a straightforward explanation for different empirical outcomes.

- Communicate findings inclusively enough for readerships from diverse backgrounds. Provide explicit definitions (e.g. for terms like “**prosody**”, “dysarthria”, or “anthropomorphism”), avoid technical jargon including abbreviations unfamiliar to other fields (e.g. synthesis algorithms, machine learning approaches, or acoustic measures), adopt scientific standards from other fields where appropriate, and discuss findings against the wider interdisciplinary literature (i.e. linking insights into pathological voices to synthetic ones and vice versa).
- Quantify naturalness whenever it could have important implications for the ecological validity of the stimulus material, even when naturalness is not the primary focus of the study. This is especially important when using acoustic manipulations which could have unintended side effects on perceived naturalness [74,116].

### **Glossary:**

**Acoustic cues:** physical and measurable features of sounds (such as voices); these may include fundamental frequency, intensity, a range of timbre cues, or temporal characteristics. Used by listeners to inform manifold impressions about voices, such as emotion, identity, age, gender or naturalness.

**Anthropomorphism:** the attribution of human characteristics, emotions, or behaviors to non-human entities

**ChatGPT:** a chatbot developed by OpenAI, based on a large language model, that generates text based on input-prompts (GPT stands for generative pre-trained transformer)

**Deepfakes:** digitally manipulated media, such as images, videos, or voice recordings, created using deep learning techniques with the goal to convincingly display the appearance of a specific individual.

1           **Deviation-based naturalness:** Conceptualization as the deviation from a reference that represents  
2           maximum naturalness.  
3  
4

5           **Dysarthria:** impairments of speech motor subsystems due to various neurological conditions such as  
6  
7           Parkinson's disease, amyotrophic lateral sclerosis (ALS), developmental conditions, strokes, or  
8  
9           traumatic brain injury.  
10  
11

12           **Human-likeness-based naturalness:** Conceptualization of naturalness by its resemblance to a real  
13           human voice.  
14  
15

16           **Laryngectomy:** surgical removal of the larynx, typically in the context of larynx cancer treatment  
17  
18

19           **Prosody:** Dynamic voice intonation, as expressed in pitch, loudness, timbre, and rhythm. Sometimes  
20  
21           also referred to as voice melody.  
22  
23

24           **Tracheoesophageal speech:** a method of vocalization following total laryngectomy via a  
25  
26           tracheoesophageal prosthesis that enables speech through esophageal vibrations.  
27  
28

29           **Uncanny valley:** a sudden feeling of eeriness evoked by humanoid robots that almost approach, but  
30  
31           do not entirely reach a human-like appearance  
32  
33

34           **Voice synthesis:** creation of computer-generated voices. Common methods are articulatory  
35  
36           synthesis, concatenative synthesis, and statistical parametric synthesis, including deep learning  
37  
38           algorithms  
39  
40  
41

#### 42           **Acknowledgements and Funding** 43 44

45  
46  
47  
48           We thank Simone Dahmen and Fatma Bilem for their support with the literature analysis, and the  
49  
50           members of the Jena Voice Research Unit (<https://www.voice.uni-jena.de/>) for helpful suggestions  
51  
52           on this project.  
53  
54  
55

The authors gratefully acknowledge the award of funding through an EU-MSCA doctoral  
1  
2 network “Voice Communication Sciences” (action 101168998, <https://www.vocs.eu.com/>).  
3  
4

CN: I dedicate this work to our stillborn son. Thanks for changing our lives.  
5  
6  
7  
8  
9  
10

11 **Declaration of interests:**  
12  
13

14 The authors declare no competing interests.  
15  
16  
17  
18  
19  
20  
21  
22

23 **References**  
24

- 25 1. Román, S. et al. (2017) The importance of food naturalness for consumers: Results of a  
26 systematic review. *Trends in Food Science & Technology* 67, 44–57. DOI:  
27 10.1016/j.tifs.2017.06.010
- 28 2. Meier, B.P. et al. (2019) Naturally better? A review of the natural-is-better bias. *Social &*  
29 *Personality Psych* 13 (8). DOI: 10.1111/spc3.12494
- 30 3. Ode, A. et al. (2009) Indicators of perceived naturalness as drivers of landscape preference.  
31 *Journal of environmental management* 90, 375–383. DOI: 10.1016/j.jenvman.2007.10.013
- 32 4. Young, A.W. et al. (2020) Face and voice perception: Understanding commonalities and  
33 differences. *Trends Cogn Sci* 24, 398–410. DOI: 10.1016/j.tics.2020.02.001
- 34 5. Rodero, E. and Lucas, I. (2023) Synthetic versus human voices in audiobooks: The human  
35 emotional intimacy effect. *New Media & Society* 25, 1746–1764. DOI:  
36 10.1177/14614448211024142
- 37 6. Rodero, E. (2017) Effectiveness, attention, and recall of human and artificial voices in an  
38 advertising story. Prosody influence and functions of voices. *Computers in Human Behavior* 77,  
39 336–346. DOI: 10.1016/j.chb.2017.08.044
- 40 7. Roswandowicz, C. et al. (2024) Cortical-striatal brain network distinguishes deepfake from real  
41 speaker identity. *Communications biology* 7, 711. DOI: 10.1038/s42003-024-06372-6
- 42 8. Lavan, N. et al. (2024) The time course of person perception from voices in the brain. *Proc Natl*  
43 *Acad Sci U S A* 121, e2318361121. DOI: 10.1073/pnas.2318361121
- 44 9. Lavan, N. (2023) How do we describe other people from voices and faces? *Cognition* 230,  
45 105253. DOI: 10.1016/j.cognition.2022.105253
- 46 10. Jiang, Z. et al. (2024) Comparison of face-based and voice-based first impressions in a Chinese  
47 sample. *Br. J. Psychol.* 115, 20–39. DOI: 10.1111/bjop.12675
- 48 11. Kühne, K. et al. (2020) The Human Takes It All: Humanlike Synthesized Voices Are Perceived as  
49 Less Eerie and More Likable. Evidence From a Subjective Ratings Study. *Frontiers in*  
50 *NeuroRobotics* 14, 1–16. DOI: 10.3389/fnbot.2020.593732
- 51 12. Ilves, M. and Surakka, V. (2013) Subjective responses to synthesised speech with lexical  
52 emotional content: the effect of the naturalness of the synthetic voice. *Behaviour & Information*  
53 *Technology* 32, 117–131. DOI: 10.1080/0144929X.2012.702285

13. Ilves, M. et al. (2011) The Effects of Emotionally Worded Synthesized Speech on the Ratings of  
1 Emotions and Voice Quality. In , pp. 588–598, Springer, Berlin, Heidelberg  
2
14. Anand, S. and Stepp, C.E. (2015) Listener Perception of Monopitch, Naturalness, and  
3 Intelligibility for Speakers With Parkinson's Disease. *J Speech Lang Hear Res* 58, 1134–1144. DOI:  
4 10.1044/2015\_JSLHR-S-14-0243  
5
- 6 Moya-Galé, G. and Levy, E.S. (2019) Parkinson's disease-associated dysarthria: prevalence,  
7 impact and management strategies. *JPRLS Volume* 9, 9–16. DOI: 10.2147/JPRLS.S168090  
8
- 9 Damico, J.S. and Ball, M.J., eds (2019) *The SAGE Encyclopedia of Human Communication Sciences*  
10 and *Disorders*, SAGE Publications, Inc  
11
- 12 Klopfenstein, M. et al. (2020) The study of speech naturalness in communication disorders: A  
13 systematic review of the literature. *Clinical Linguistics & Phonetics* 34, 327–338. DOI:  
14 10.1080/02699206.2019.1652692  
15
- 16 Frankford, S.A. et al. (2024) Contributions of Speech Timing and Articulatory Precision to Listener  
17 Perceptions of Intelligibility and Naturalness in Parkinson's Disease. *J Speech Lang Hear Res* 67,  
18 2951–2963. DOI: 10.1044/2024\_JSLHR-23-00802  
19
- 20 Euler, H.A. et al. (2021) Speech restructuring group treatment for 6-to-9-year-old children who  
21 stutter: A therapeutic trial. *Journal of communication disorders* 89, 106073. DOI:  
22 10.1016/j.jcomdis.2020.106073  
23
- 24 Hardy, T.L.D. et al. (2020) Acoustic Predictors of Gender Attribution, Masculinity-Femininity, and  
25 Vocal Naturalness Ratings Amongst Transgender and Cisgender Speakers. *Journal of Voice* 34,  
26 300.e11-300.e26. DOI: 10.1016/j.jvoice.2018.10.002  
27
- 28 Hyppa-Martin, J. et al. (2024) A large-scale comparison of two voice synthesis techniques on  
29 intelligibility, naturalness, preferences, and attitudes toward voices banked by individuals with  
30 amyotrophic lateral sclerosis. *Augmentative and Alternative Communication* 40, 31–45. DOI:  
31 10.1080/07434618.2023.2262032  
32
- 33 Nass, C. et al. (1994) Computers are social actors. In *Proceedings of the SIGCHI conference on*  
34 *Human factors in computing systems celebrating interdependence - CHI '94*, ACM Press  
35
- 36 Seaborn, K. et al. (2021) Voice in Human-Agent Interaction. *ACM Comput. Surv.* 54, 1–43. DOI:  
37 10.1145/3386867  
38
- 39 Triantafyllopoulos, A. et al. (2023) An overview of affective speech synthesis and conversion in  
40 the deep learning era. *Proceedings of the IEEE*, 1355–1381  
41
- 42 Schreibelmayr, S. and Mara, M. (2022) Robot Voices in Daily Life: Vocal Human-Likeness and  
43 Application Context as Determinants of User Acceptance. *Frontiers in Psychology* 13, 1–17. DOI:  
44 10.3389/fpsyg.2022.787499  
45
- 46 Baird, A. et al. (2018) The Perception and Analysis of the Likeability and Human Likeness of  
47 Synthesized Speech. In *Interspeech 2018*, pp. 2863–2867, ISCA  
48
- 49 Lee, E.-J. (2010) The more humanlike, the better? How speech type and users' cognitive style  
50 affect social responses to computers. *Computers in Human Behavior* 26, 665–672. DOI:  
51 10.1016/j.chb.2010.01.003  
52
- 53 Lu, L. et al. (2021) Leveraging "human-likeness" of robotic service at restaurants. *International*  
54 *Journal of Hospitality Management* 94, 1–9. DOI: 10.1016/j.ijhm.2020.102823  
55
- 56 Cambre, J. and Kulkarni, C. (2019) One Voice Fits All? *Proc. ACM Hum.-Comput. Interact.* 3, 1–19.  
57 DOI: 10.1145/3359325  
58
- 59 Eyssel, F. et al. (2012) 'If you sound like me, you must be more human'. In *HRI' 12. Proceedings of*  
60 *the seventh annual ACM/IEEE Conference on Human-Robot Interaction : March 5-8, 2012 Boston,*  
61 *Massachusetts, USA* (Yanco, H. et al., eds), pp. 125–126, Association for Computing Machinery  
62
- 63
- 64
- 65

- 1           31. Im, H. et al. (2023) Let voice assistants sound like a machine: Voice and task type effects on  
2           perceived fluency, competence, and consumer attitude. *Computers in Human Behavior* 145,  
3           107791. DOI: 10.1016/j.chb.2023.107791
- 4           32. McGinn, C. and Torre, I. (2019 - 2019) Can you Tell the Robot by the Voice? An Exploratory Study  
5           on the Role of Voice in the Perception of Robots. In *2019 14th ACM/IEEE International  
6           Conference on Human-Robot Interaction (HRI)*, pp. 211–221, IEEE
- 7           33. Mitchell, W.J. et al. (2011) A mismatch in the human realism of face and voice produces an  
8           uncanny valley. *i-Perception* 2, 10–12. DOI: 10.1068/i0415
- 9           34. Yorkston, K.M. et al. (1999) *Management of motor speech disorders in children and adults*, Pro-  
10           ed Austin, TX
- 11           35. Mawalim, C.O. et al. (2022) Speaker anonymization by modifying fundamental frequency and x-  
12           vector singular value. *Computer Speech & Language* 73, 1–17. DOI: 10.1016/j.csl.2021.101326
- 13           36. Hu, P. et al. (2021) Dual humanness and trust in conversational AI: A person-centered approach.  
14           *Computers in Human Behavior* 119, 106727. DOI: 10.1016/j.chb.2021.106727
- 15           37. Nusbaum, H.C. et al. (1997) Measuring the naturalness of synthetic speech. *International Journal  
16           of Speech Technology* 2, 7–19
- 17           38. Mayo, C. et al. (2011) Listeners' weighting of acoustic cues to synthetic speech naturalness: A  
18           multidimensional scaling analysis. *Speech Commun* 53, 311–326. DOI:  
19           10.1016/j.specom.2010.10.003
- 20           39. Abdulrahman, A. and Richards, D. (2022) Is Natural Necessary? Human Voice versus Synthetic  
21           Voice for Intelligent Virtual Agents. *MTI* 6, 51. DOI: 10.3390/mti6070051
- 22           40. Urakami, J. et al. (2020) The Effect of Naturalness of Voice and Empathic Responses on  
23           Enjoyment, Attitudes and Motivation for Interacting with a Voice User Interface. In *Human-  
24           Computer Interaction. Multimodal and Natural Interaction* (Kurosu, M., ed), pp. 244–259,  
25           Springer International Publishing
- 26           41. Velner, E. et al. (2020) Intonation in Robot Speech. In *Proceedings of the 2020 ACM/IEEE  
27           International Conference on Human-Robot Interaction* (Belpaeme, T. et al., eds), pp. 569–578,  
28           ACM
- 29           42. Yamasaki, R. et al. (2017) Perturbation Measurements on the Degree of Naturalness of  
30           Synthesized Vowels. *Journal of Voice* 31, 389.e1-389.e8. DOI: 10.1016/j.jvoice.2016.09.020
- 31           43. Ko, S. et al. (2023) The Effects of Robot Voices and Appearances on Users' Emotion Recognition  
32           and Subjective Perception. *Int. J. Human. Robot.* 20. DOI: 10.1142/S0219843623500019
- 33           44. Abur, D. et al. (2021) Feedback and Feedforward Auditory-Motor Processes for Voice and  
34           Articulation in Parkinson's Disease. *J Speech Lang Hear Res* 64, 4682–4694. DOI:  
35           10.1044/2021\_JSLHR-21-00153
- 36           45. Klopfenstein, M. (2015) Relationship between acoustic measures and speech naturalness ratings  
37           in Parkinson's disease: A within-speaker approach. *Clinical Linguistics & Phonetics* 29, 938–954.  
38           DOI: 10.3109/02699206.2015.1081293
- 39           46. Klopfenstein, M. (2016) Speech naturalness ratings and perceptual correlates of highly natural  
40           and unnatural speech in hypokinetic dysarthria secondary to Parkinson's disease. *JRCD* 7, 123–  
41           146. DOI: 10.1558/jircd.v7i1.27932
- 42           47. Moya-Galé, G. et al. (2024) Perceptual consequences of online group speech treatment for  
43           individuals with Parkinson's disease: A pilot study case series. *International Journal of Speech-  
44           Language Pathology*, 1–16. DOI: 10.1080/17549507.2024.2330538
- 45           48. Eadie, T.L. and Doyle, P.C. (2002) Direct Magnitude Estimation and Interval Scaling of  
46           Naturalness and Severity in Tracheoesophageal (TE) Speakers. *J Speech Lang Hear Res* 45, 1088–  
47           1096. DOI: 10.1044/1092-4388(2002/087)

- 1           49. Eadie, T.L. et al. (2008) Influence of speaker gender on listener judgments of tracheoesophageal  
2           speech. *Journal of Voice* 22, 43–57. DOI: 10.1016/j.jvoice.2006.08.008  
3           50. Yorkston, K.M. et al. (1990) The effect of rate control on the intelligibility and naturalness of  
4           dysarthric speech. *The Journal of speech and hearing disorders* 55, 550–560. DOI:  
5           10.1044/jshd.5503.550  
6           51. Schölderle, T. et al. (2023) Speech Naturalness in the Assessment of Childhood Dysarthria.  
7           *American Journal of Speech-language Pathology* 32, 1633–1643. DOI: 10.1044/2023\_AJSLP-23-  
8           00023  
9           52. Lehner, K. and Ziegler, W. (2022) Clinical measures of communication limitations in dysarthria  
10          assessed through crowdsourcing: specificity, sensitivity, and retest-reliability. *Clinical Linguistics  
11          & Phonetics* 36, 988–1009. DOI: 10.1080/02699206.2021.1979658  
12          53. Vogel, A.P. et al. (2019) Speech treatment improves dysarthria in multisystemic ataxia: a rater-  
13          blinded, controlled pilot-study in ARSACS. *Journal of neurology* 266, 1260–1266. DOI:  
14          10.1007/s00415-019-09258-4  
15          54. Jones, H.N. et al. (2019) Auditory-Perceptual Speech Features in Children With Down Syndrome.  
16          *American journal on intellectual and developmental disabilities* 124, 324–338. DOI:  
17          10.1352/1944-7558-124.4.324  
18          55. Assmann, P.F. et al. (2006) Effects of frequency shifts on perceived naturalness and gender  
19          information in speech. In *INTERSPEECH*  
20          56. Venkatraman, A. and Sivasankar, M.P. (2018) Continuous Vocal Fry Simulated in Laboratory  
21          Subjects: A Preliminary Report on Voice Production and Listener Ratings. *American Journal of  
22          Speech-language Pathology* 27, 1539–1545. DOI: 10.1044/2018\_AJSLP-17-0212  
23          57. Kapolowicz, M.R. et al. (2022) Effects of Spectral Envelope and Fundamental Frequency Shifts on  
24          the Perception of Foreign-Accented Speech. *Language and speech* 65, 418–443. DOI:  
25          10.1177/00238309211029679  
26          58. Tamagawa, R. et al. (2011) The Effects of Synthesized Voice Accents on User Perceptions of  
27          Robots. *Int J of Soc Robotics* 3, 253–262. DOI: 10.1007/s12369-011-0100-4  
28          59. Mackey, L.S. et al. (1997) Effect of speech dialect on speech naturalness ratings: a systematic  
29          replication of Martin, Haroldson, and Triden (1984). *J Speech Lang Hear Res* 40, 349–360. DOI:  
30          10.1044/jslhr.4002.349  
31          60. Goy, H. et al. (2016) Effects of age on speech and voice quality ratings. *The Journal of the  
32          Acoustical Society of America* 139, 1648. DOI: 10.1121/1.4945094  
33          61. Coughlin-Woods, S. et al. (2005) Ratings of speech naturalness of children ages 8–16 years.  
34          *Percept Motor Skill* 100, 295–304. DOI: 10.2466/pms.100.2.295-304  
35          62. Baird, A. et al. (2017) Perception of Paralinguistic Traits in Synthesized Voices. In *Proceedings of  
36          the 12th International Audio Mostly Conference on Augmented and Participatory Sound and  
37          Music Experiences* (Fazekas, G. et al., eds), pp. 1–5, ACM  
38          63. Merritt, B. and Bent, T. (2020) Perceptual Evaluation of Speech Naturalness in Speakers of  
39          Varying Gender Identities. *J Speech Lang Hear Res* 63, 2054–2069. DOI: 10.1044/2020\_JSLHR-19-  
40          00337  
41          64. Baird, A. et al. (2018) The Perception of Vocal Traits in Synthesized Voices: Age, Gender, and  
42          Human Likeness. *J. Audio Eng. Soc.* 66, 277–285. DOI: 10.17743/jaes.2018.0023  
43          65. Aylett, M.P. et al. (2020) Speech Synthesis for the Generation of Artificial Personality. *IEEE Trans.  
44          Affective Comput.* 11, 361–372. DOI: 10.1109/TAFFC.2017.2763134  
45          66. Kramer, R.S.S. et al. (2024) The psychometrics of rating facial attractiveness using different  
46          response scales. *Perception* 53, 645–660. DOI: 10.1177/03010066241256221  
47          67. Martin, R.R. et al. (1984) Stuttering and speech naturalness. *The Journal of speech and hearing  
48          disorders* 49, 53–58. DOI: 10.1044/jshd.4901.53

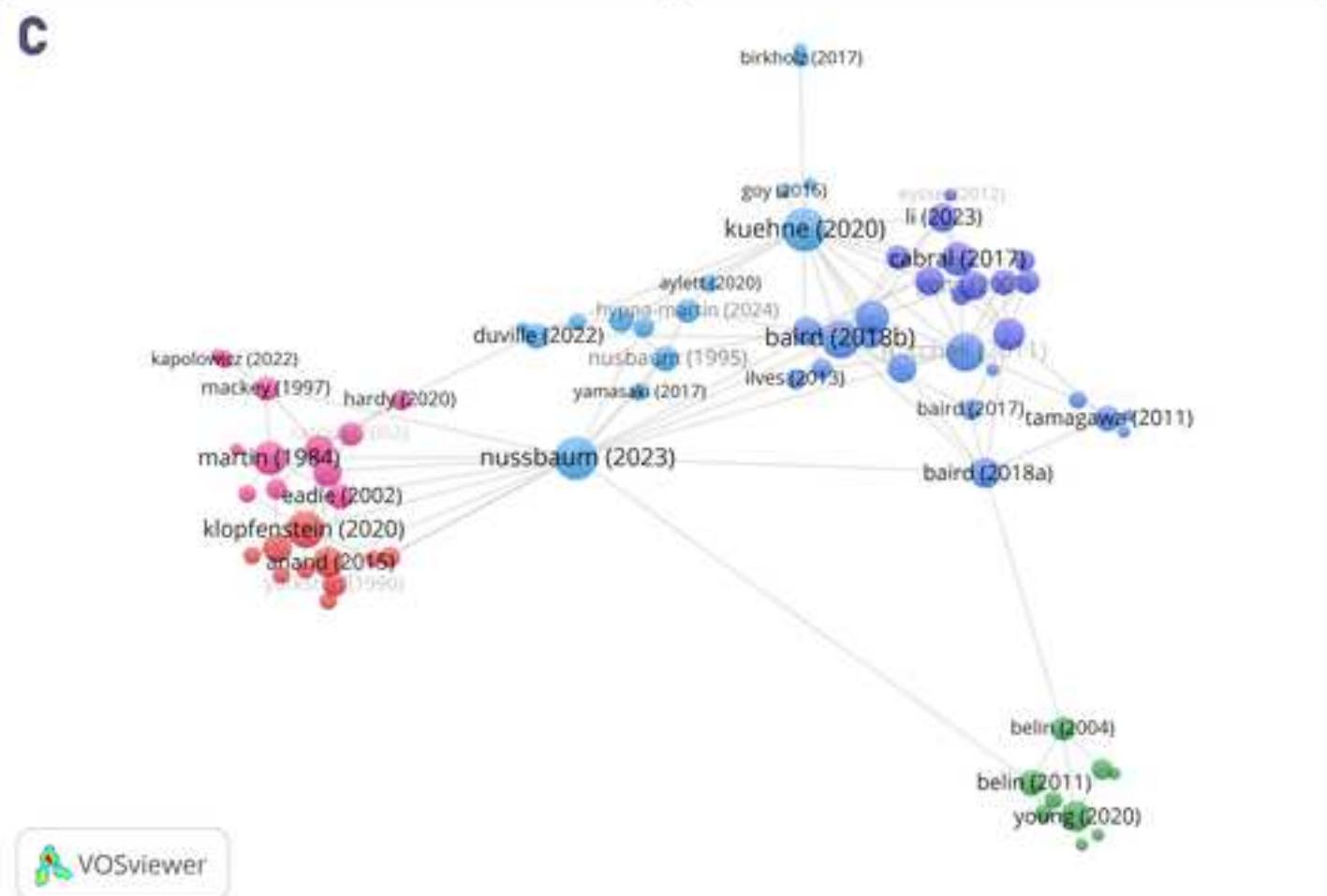
- 1           68. van Eck, N.J. and Waltman, L. (2010) Software survey: VOSviewer, a computer program for  
2           bibliometric mapping. *Scientometrics* 84, 523–538. DOI: 10.1007/s11192-009-0146-3  
3           69. van der Linden, S. (2023) *Foolproof: Why we fall for misinformation and how to build immunity*,  
4           WW Norton & Company.  
5           70. Fiske, S.T. (2018) Stereotype Content: Warmth and Competence Endure. *Curr Dir Psychol Sci* 27,  
6           67–73. DOI: 10.1177/0963721417738825  
7           71. Todorov, A. et al. (2008) Understanding evaluation of faces on social dimensions. *Trends Cogn Sci*  
8           12, 455–460. DOI: 10.1016/j.tics.2008.10.001  
9           72. Sutherland, C.A.M. et al. (2013) Social inferences from faces: ambient images generate a three-  
10          dimensional model. *Cognition* 127, 105–118. DOI: 10.1016/j.cognition.2012.12.001  
11          73. Sutherland, C.A.M. et al. (2016) Integrating social and facial models of person perception:  
12          Converging and diverging dimensions. *Cognition* 157, 257–267. DOI:  
13          10.1016/j.cognition.2016.09.006  
14          74. Nussbaum, C. et al. (2023) Perceived naturalness of emotional voice morphs. *Cognition &*  
15          *Emotion*, 1–17. DOI: 10.1080/02699931.2023.2200920  
16          75. Mori, M. et al. (2012) The Uncanny Valley. *IEEE Robot. Automat. Mag.* 19, 98–100. DOI:  
17          10.1109/mra.2012.2192811  
18          76. Romportl, J. (2014) Speech Synthesis and Uncanny Valley. In *Text, speech and dialogue* (Horák, A.  
19          et al., eds), pp. 595–602, Springer International Publishing  
20          77. Diel, A. and Lewis, M. (2024) Deviation from typical organic voices best explains a vocal uncanny  
21          valley. *Computers in Human Behavior Reports* 14, 100430. DOI: 10.1016/j.chbr.2024.100430  
22          78. van Prooije, T. et al. (2024) Perceptual and Acoustic Analysis of Speech in Spinocerebellar ataxia  
23          Type 1. *Cerebellum*, 112–120. DOI: 10.1007/s12311-023-01513-9  
24          79. Moore, B.C.J. and Tan, C.-T. (2003) Perceived naturalness of spectrally distorted speech and  
25          music. *The Journal of the Acoustical Society of America* 114, 408–419. DOI: 10.1121/1.1577552  
26          80. Rao M V, A. et al. (2018) Effect of source filter interaction on isolated vowel-consonant-vowel  
27          perception. *The Journal of the Acoustical Society of America* 144, EL95. DOI: 10.1121/1.5049510  
28          81. Ratcliff, A. et al. (2002) Factors influencing ratings of speech naturalness in augmentative and  
29          alternative communication. *Augmentative and Alternative Communication* 18, 11–19. DOI:  
30          10.1080/aac.18.1.11.19  
31          82. Meltzner, G.S. and Hillman, R.E. (2005) Impact of Aberrant Acoustic Properties on the Perception  
32          of Sound Quality in Electrolarynx Speech. *J Speech Lang Hear Res* 48, 766–779. DOI:  
33          10.1044/1092-4388(2005/053)  
34          83. Andics, A. et al. (2010) Neural mechanisms for voice recognition. *Neuroimage* 52, 1528–1540.  
35          DOI: 10.1016/j.neuroimage.2010.05.048  
36          84. Valentine, T. et al. (2016) Face-space: A unifying concept in face recognition research. *Q J Exp*  
37          *Psychol (Hove)* 69, 1996–2019. DOI: 10.1080/17470218.2014.990392  
38          85. Lima, C.F. et al. (2021) Authentic and posed emotional vocalizations trigger distinct facial  
39          responses. *Cortex* 141, 280–292. DOI: 10.1016/j.cortex.2021.04.015  
40          86. Sarzedas, J. et al. (2024) Blindness influences emotional authenticity perception in voices:  
41          Behavioral and ERP evidence. *Cortex* 172, 254–270. DOI: 10.1016/j.cortex.2023.11.005  
42          87. Anikin, A. and Lima, C.F. (2017) Perceptual and acoustic differences between authentic and  
43          acted nonverbal emotional vocalizations. *Q J Exp Psychol (Hove)* 71, 622–641. DOI:  
44          10.1080/17470218.2016.1270976  
45          88. Kachel, S. et al. (2020) Gender (Conformity) Matters: Cross-Dimensional and Cross-Modal  
46          Associations in Sexual Orientation Perception. *Journal of Language and Social Psychology* 39, 40–  
47          66. DOI: 10.1177/0261927X19883902

- 1 89. Mills, M. et al. (2017) Expanding the evidence: Developments and innovations in clinical practice,  
2 training and competency within voice and communication therapy for trans and gender diverse  
3 people. *International Journal of Transgenderism* 18, 328–342. DOI:  
4 10.1080/15532739.2017.1329049
- 5 90. Eiff, C.I. von et al. (2022) Crossmodal benefits to vocal emotion perception in cochlear implant  
6 users. *iScience* 25, 105711. DOI: 10.1016/j.isci.2022.105711
- 7 91. Schweinberger, S.R. and Eiff, C.I. von (2022) Enhancing socio-emotional communication and  
8 quality of life in young cochlear implant recipients: Perspectives from parameter-specific  
9 morphing and caricaturing. *Frontiers in Neuroscience* 16, 956917. DOI:  
10 10.3389/fnins.2022.956917
- 11 92. Yamagishi, J. et al. (2012) Speech synthesis technologies for individuals with vocal disabilities:  
12 Voice banking and reconstruction. *Acoust. Sci. & Tech.* 33, 1–5. DOI: 10.1250/ast.33.1
- 13 93. Belin, P. et al. (2004) Thinking the voice: neural correlates of voice perception. *Trends Cogn Sci* 8,  
14 129–135. DOI: 10.1016/j.tics.2004.01.008
- 15 94. Belin, P. et al. (2011) Understanding voice perception. *Br. J. Psychol.* 102, 711–725. DOI:  
16 10.1111/j.2044-8295.2011.02041.x
- 17 95. Lavan, N. and McGettigan, C. (2023) A model for person perception from familiar and unfamiliar  
18 voices. *Commun Psychol* 1, 1–11. DOI: 10.1038/s44271-023-00001-4
- 19 96. Staib, M. and Frühholz, S. (2023) Distinct functional levels of human voice processing in the  
20 auditory cortex. *Cerebral Cortex* 33, 1170–1185. DOI: 10.1093/cercor/bhac128
- 21 97. Staib, M. and Frühholz, S. (2021) Cortical voice processing is grounded in elementary sound  
22 analyses for vocalization relevant sound patterns. *Progress in neurobiology* 200, 101982. DOI:  
23 10.1016/j.pneurobio.2020.101982
- 24 98. Pinheiro, A.P. et al. (2021) Emotional authenticity modulates affective and social trait inferences  
25 from voices. *Philosophical transactions of the Royal Society of London. Series B, Biological  
26 sciences* 376, 20200402. DOI: 10.1098/rstb.2020.0402
- 27 99. Duville, M.M. et al. (2022) Neuronal and behavioral affective perceptions of human and  
28 naturalness-reduced emotional prosodies. *Frontiers in computational neuroscience* 16, 1022787.  
29 DOI: 10.3389/fncom.2022.1022787
- 30 100. Duville, M.M. et al. (2024) Improved emotion differentiation under reduced acoustic variability  
31 of speech in autism. *BMC medicine* 22, 121. DOI: 10.1186/s12916-024-03341-y
- 32 101. Nussbaum, C. et al. (2022) Contributions of fundamental frequency and timbre to vocal emotion  
33 perception and their electrophysiological correlates. *Social Cognitive and Affective Neuroscience*  
34 17, 1145–1154. DOI: 10.1093/scan/nsac033
- 35 102. Kosilo, M. et al. (2021) The neural basis of authenticity recognition in laughter and crying.  
36 *Scientific reports* 11, 23750. DOI: 10.1038/s41598-021-03131-z
- 37 103. Conde, T. et al. (2022) The time course of emotional authenticity detection in nonverbal  
38 vocalizations. *Cortex; a journal devoted to the study of the nervous system and behavior* 151,  
39 116–132. DOI: 10.1016/j.cortex.2022.02.016
- 40 104. Miller, E.J. et al. (2023) How do people respond to computer-generated versus human faces? A  
41 systematic review and meta-analyses. *Computers in Human Behavior Reports*, 100283. DOI:  
42 10.1016/j.chbr.2023.100283
- 43 105. Miller, E.J. et al. (2023) AI Hyperrealism: Why AI Faces Are Perceived as More Real Than Human  
44 Ones. *Psychol Sci* 34, 1390–1403. DOI: 10.1177/09567976231207095
- 45 106. Cabral, J.P. et al. (2017) The Influence of Synthetic Voice on the Evaluation of a Virtual Character.  
46 In *Interspeech 2017*, pp. 229–233, ISCA
- 47 107. Ehret, J. et al. (2021) Do Prosody and Embodiment Influence the Perceived Naturalness of  
48 Conversational Agents' Speech? *ACM Trans. Appl. Percept.* 18, 1–15. DOI: 10.1145/3486580

- 108.Ferstl, Y. et al. (2021) Human or Robot? Investigating voice, appearance and gesture motion  
realism of conversational social agents. In *Proceedings of the 21th ACM International Conference  
on Intelligent Virtual Agents*, pp. 76–83, ACM
- 109.Gong, L. and Nass, C. (2007) When a Talking-Face Computer Agent is Half-Human and Half-  
Humanoid: Human Identity and Consistency Preference. *Human Comm Res* 33, 163–193. DOI:  
10.1111/j.1468-2958.2007.00295.x
- 110.Higgins, D. et al. (2022) Sympathy for the digital: Influence of synthetic voice on affinity, social  
presence and empathy for photorealistic virtual humans. *Computers & Graphics* 104, 116–128.  
DOI: 10.1016/j.cag.2022.03.009
- 111.Li, M. et al. (2023) Effects of robot gaze and voice human-likeness on users' subjective  
perception, visual attention, and cerebral activity in voice conversations. *Computers in Human  
Behavior* 141, 107645. DOI: 10.1016/j.chb.2022.107645
- 112.Parmar, D. et al. (2022) Designing Empathic Virtual Agents: Manipulating Animation, Voice,  
Rendering, and Empathy to Create Persuasive Agents. *Autonomous agents and multi-agent  
systems* 36, 1–24. DOI: 10.1007/s10458-021-09539-1
- 113.Sarigul, B. and Urgen, B.A. (2023) Audio–Visual Predictive Processing in the Perception of  
Humans and Robots. *Int J of Soc Robotics* 15, 855–865. DOI: 10.1007/s12369-023-00990-6
- 114.Lowry, H. et al. (2013) Behavioural responses of wildlife to urban environments. *Biological  
reviews of the Cambridge Philosophical Society* 88, 537–549. DOI: 10.1111/brv.12012
- 115.Kauk, J. et al. (2024) The adaptive community-response (ACR) method for collecting  
misinformation on social media. *J Big Data* 11, 1–32. DOI: 10.1186/s40537-024-00894-w
- 116.Malisz, Z. et al. (2020) Modern speech synthesis for phonetic sciences: a discussion and an  
evaluation, 487–491. DOI: 10.31234/osf.io/dxvhc

## Outstanding questions:

- Vocal communication is abundant in the animal kingdom, and many animals manipulate their vocal behavior in an adaptive manner – is there demand for a comparative perspective on voice naturalness?
- How is a listener's perception of naturalness shaped through experience (e.g., with voice assistants, smart home devices, or patients with voice disorders)?
- With respect to the present conceptual framework, (how) are human-likeness based naturalness and deviation-based naturalness dissociable in the brain?
- In the trade-off between precise experimental control and open field recordings, can we identify converging evidence for how and when reduced naturalness in voices critically affects the ecological validity of research? In depth, will we need a dynamic definition of ecological validity in view of an ever more digital world of social interaction?
- Are natural voices always preferred, or is naturalness preference context dependent? Can natural voices impede rather than promote communication success in some situations?
- Many domains of social perception are characterized by individual variability, but it is unclear whether there are substantial individual differences in the tolerance of or preference for unnatural voice features. If so, can these be related to other domains of auditory cognition, or to other person traits?
- To what extent is naturalness perception affected by factors such as age, gender, or cultural background?



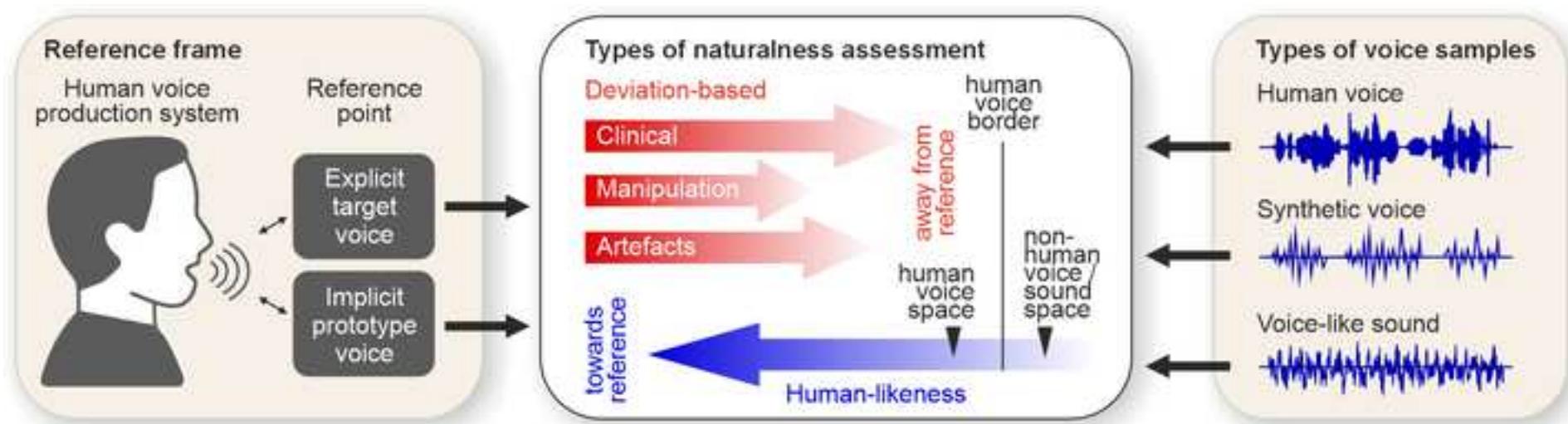
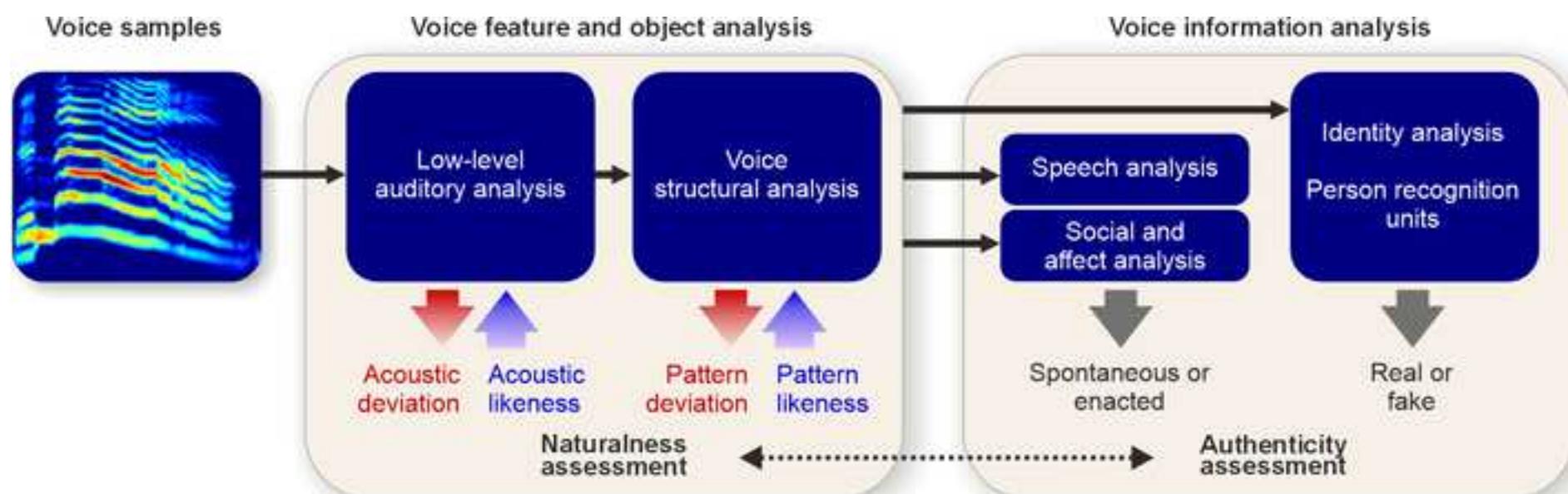


Figure 3

[Click here to access/download;Figure;Figure3.tif](#)

## CELL PRESS DECLARATION OF INTERESTS POLICY

Transparency is essential for a reader's trust in the scientific process and for the credibility of published articles. At Cell Press, we feel that disclosure of competing interests is a critical aspect of transparency. Therefore, we require a "declaration of interests" section in which all authors disclose any financial or other interests related to the submitted work that (1) could affect or have the perception of affecting the author's objectivity or (2) could influence or have the perception of influencing the content of the article.

### **What types of articles does this apply to?**

We require that you disclose competing interests for all submitted content by completing and submitting the form below. We also require that you include a "declaration of interests" section in the text of all articles even if there are no interests to declare.

### **What should I disclose?**

We require that you and all authors disclose any personal financial interests (e.g., stocks or shares in companies with interests related to the submitted work or consulting fees from companies that could have interests related to the work), professional affiliations, advisory positions, board memberships (including membership on a journal's advisory board when publishing in that journal), or patent applications and/or registrations that are related to the subject matter of the contribution. As a guideline, you need to declare an interest for (1) any affiliation associated with a payment or financial benefit exceeding \$10,000 p.a. or 5% ownership of a company or (2) research funding by a company with related interests. You do not need to disclose diversified mutual funds, 401ks, or investment trusts.

Authors should also disclose relevant financial interests of immediate family members. Cell Press uses the Public Health Service definition of "immediate family member," which includes spouse and dependent children.

### **Where do I declare competing interests?**

Competing interests should be disclosed on this form as well as in a "declaration of interests" section in the manuscript. This section should include financial or other competing interests as well as affiliations that are not included in the author list. Examples of "declaration of interests" language include:

"AUTHOR is an employee and shareholder of COMPANY."

"AUTHOR is a founder of COMPANY and a member of its scientific advisory board."

**NOTE:** Primary affiliations should be included with the author list and do not need to be included in the "declaration of interests" section. Funding sources should be included in the "acknowledgments" section and also do not need to be included in the "declaration of interests" section. (A small number of front-matter article types do not include an "acknowledgments" section. For these articles, reporting of funding sources is not required.)

### **What if there are no competing interests to declare?**

If you have no competing interests to declare, please note that in the "declaration of interests" section with the following wording:

"The authors declare no competing interests."

## CELL PRESS DECLARATION OF INTERESTS FORM

If submitting materials via Editorial Manager, please complete this form and upload with your initial submission. Otherwise, please email as an attachment to the editor handling your manuscript.

***Please complete each section of the form and insert any necessary “declaration of interests” statement in the text box at the end of the form. A matching statement should be included in a “declaration of interests” section in the manuscript.***

### **Institutional affiliations**

We require that you list the current institutional affiliations of all authors, including academic, corporate, and industrial, on the title page of the manuscript. ***Please select one of the following:***

- All affiliations are listed on the title page of the manuscript.
- I or other authors have additional affiliations that we have noted in the “declaration of interests” section of the manuscript and on this form below.

### **Funding sources**

We require that you disclose all funding sources for the research described in this work. ***Please confirm the following:***

- All funding sources for this study are listed in the “acknowledgments” section of the manuscript.\*

\*A small number of front-matter article types do not include an “acknowledgments” section. For these, reporting funding sources is not required.

### **Competing financial interests**

We require that authors disclose any financial interests and any such interests of immediate family members, including financial holdings, professional affiliations, advisory positions, board memberships, receipt of consulting fees, etc., that:

- (1) could affect or have the perception of affecting the author’s objectivity, or
- (2) could influence or have the perception of influencing the content of the article.

***Please select one of the following:***

- We, the authors and our immediate family members, have no financial interests to declare.
- We, the authors, have noted any financial interests in the “declaration of interests” section of the manuscript and on this form below, and we have noted interests of our immediate family members.

**Advisory/management and consulting positions**

We require that authors disclose any position, be it a member of a board or advisory committee or a paid consultant, that they have been involved with that is related to this study. We also require that members of our journal advisory boards disclose their position when publishing in that journal. **Please select one of the following:**

- We, the authors and our immediate family members, have no positions to declare and are not members of the journal's advisory board.
- The authors and/or their immediate family members have management/advisory or consulting relationships noted in the "declaration of interests" section of the manuscript and on this form below.

**Patents**

We require that you disclose any patent applications and/or registrations related to this work by any of the authors or their institutions. **Please select one of the following:**

- We, the authors and our immediate family members, have no related patent applications or registrations to declare.
- We, the authors, have a patent application and/or registration related to this work, which is noted in the "declaration of interests" section of the manuscript and on this form below, and we have noted the patents of immediate family members.

***Please insert any "declaration of interests" statements in this space.*** This exact text should also be included in the "declaration of interests" section of the manuscript. If no authors have a competing interest, please insert the text, "The authors declare no competing interests."

The authors declare no competing interests.

- On behalf of all authors, I declare that I have disclosed all competing interests related to this work. If any exist, they have been included in the "declaration of interests" section of the manuscript.**