

## Review

## Understanding voice naturalness

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

The 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 (i) conceptual underspecification, (ii) heterogeneous operationalization, (iii) lack of exchange between research on human and synthetic voices, and (iv) 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 in both empirical findings and theoretical models. Finally, it identifies gaps in current understanding of voice naturalness and sketches perspectives for empirical progress.

**Naturalness: a prominent aspect of voice perception**

Naturalness plays a significant role in how we perceive our environment through sight, sounds, smell, 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 artificial intelligence (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 in 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 prefer a more natural-sounding voice, even at the cost of some loss in intelligibility, both as users and as 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

**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 cross-refer neither to one another nor to current voice perception theory.

This review 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.

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

<sup>4</sup>Department of Psychology, University of Oslo, 0371 Oslo, Norway

<sup>5</sup>Cognitive and Affective Neuroscience Unit, University of Zurich, 8050 Zurich, Switzerland

<sup>6</sup>Swiss Center for Affective Sciences, University of Geneva, 1222 Geneva, Switzerland

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

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].

\*Correspondence:  
[christine.nussbaum@uni-jena.de](mailto:christine.nussbaum@uni-jena.de)  
(C. Nussbaum).

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: (i) conceptual underspecification, (ii) heterogeneous operationalization, (iii) lack of exchange between research domains, and (iv) 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

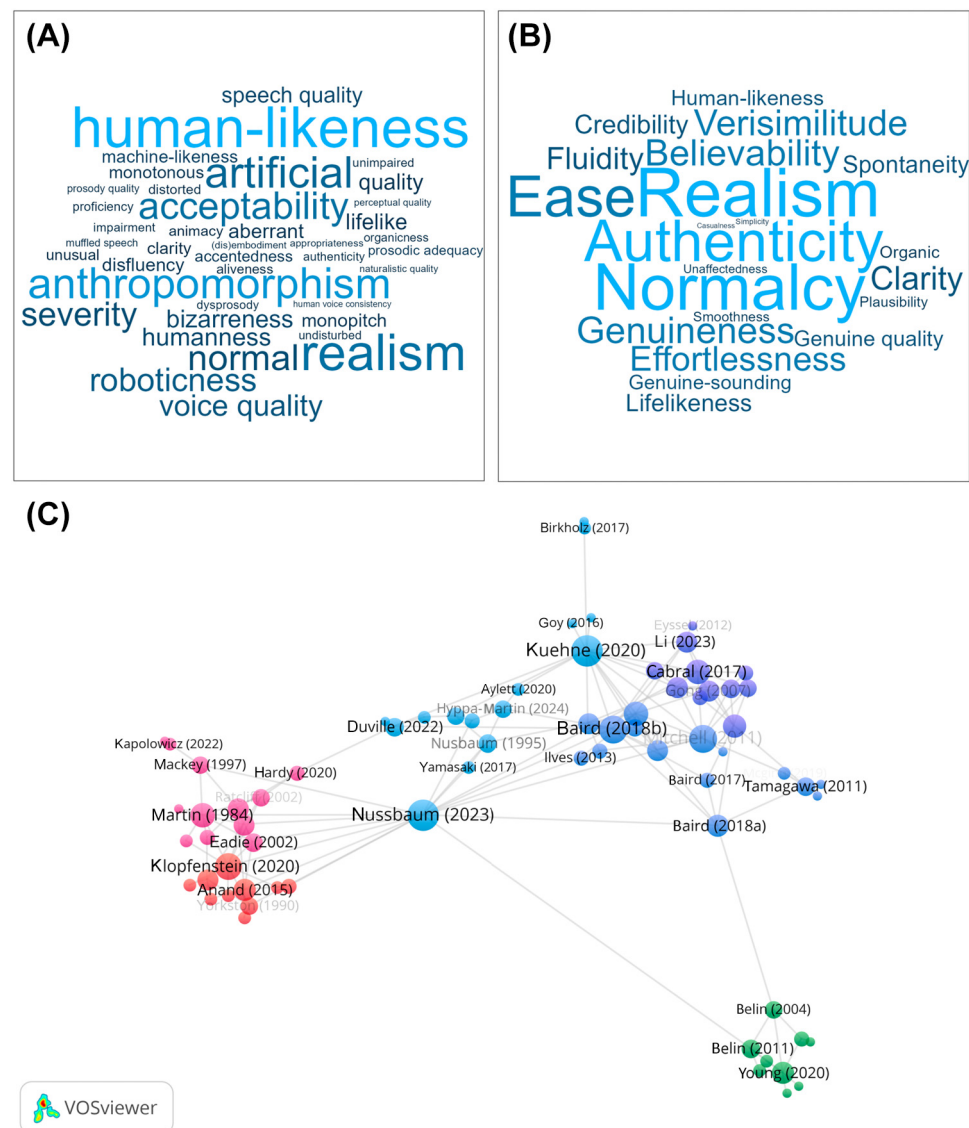
### Conceptual underspecification

Voice naturalness lacks a consistent definition and terminology in the literature (Figure 1A,B). Many papers do not even provide an explicit definition of naturalness (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 (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’ [34]. By 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, voice naturalness is a perceptual and subjective measure [37]. Second, 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 versus 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



**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 in the main text). Word size represents number of occurrences. (B) A similar word cloud but generated by ChatGPT (<https://chatgpt.com/?oai>), when prompted to generate ten 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 the associated OSF repository (<https://doi.org/10.17605/OSF.IO/ASFQV>). (C) A bibliographic network visualization using VOSviewer [68], covering publications related to voice naturalness across different domains and ten basic voice theory papers (see the supplemental information online for a full list of references). Each colored dot represents a publication and gray 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 the associated OSF repository (<https://doi.org/10.17605/OSF.IO/ASFQV>).

## 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 on the basis of input prompts. (GPT stands for generative pretrained 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.

**Deviation-based naturalness:** conceptualization as the deviation from a reference that represents maximum naturalness.

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

**Human-likeness-based naturalness:** conceptualization of naturalness by its resemblance to a real human voice.

**Laryngectomy:** surgical removal of the larynx, typically in the context of laryngeal cancer treatment.

**Prosody:** dynamic voice intonation, as expressed in pitch, loudness, timbre, and rhythm. Sometimes also referred to as 'voice melody.'

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

**Uncanny valley:** a sudden feeling of eeriness evoked by humanoid robots that almost approach but do not entirely reach a human-like appearance.

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

### Box 1. A field in numbers

For a more systematic overview of scientific insights into naturalness in voices, a focused literature search on Web of Science was conducted on April 26, 2023, using the search terms 'naturalness AND voice' or 'human-likeness AND voice,' which was repeated on May 28, 2024, to detect the most recent papers. This initial search resulted in 339 articles, to which the following inclusion criteria were applied: (i) language of publication was English; (ii) papers were published in peer-reviewed journals or as a conference contribution; (iii) voice naturalness/human-likeness was either measured or manipulated; (iv) 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; and (v) finally, the search was focused on spoken utterances, excluding singing voices and nonlinguistic 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 the supplemental information online and the associated OSF repository (<https://doi.org/10.17605/OSF.IO/ASFQV>).

In total, 72 articles were identified, covering a time range from 1984 to 2024. Thirty-eight (53%) were published in the past 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, six human-manipulated, and five 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 in the main text for examples and the supplemental information online 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 in the main text. Subsequently, these were compared with 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.

Table 1. Example definitions of deviation-based and human-likeness-based voice naturalness<sup>a</sup>

Conceptualization	Definition	Refs
Deviation-based naturalness	'Naturalness was 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.' (p. 4687)	[44]
	'Speech naturalness can be described as how the speech of a person with a speech disorder compares with that of typical speech or, in the case of an acquired disorder, how an individual's speech compares to its premorbid state' (p. 1134)	[14]
	'Speech naturalness refers to a rather broad perceptual impression representing the overall quality of a person's speech output in relation to what is conceptualized as normal or natural' (p. 1633/1634)	[51]
	'[...] degree to which individuals sound 'different' from healthy peers' (p. 1265)	[53]
Human-likeness-based naturalness	'Human likeness has been used [...] to describe how accurately the machine is able to imitate a human.' (p. 2864)	[26]
	'Naturalness refers to whether synthetic speech is perceived as uniquely human, despite being computer-generated.' (p. 5)	[21]
	'Natural speech is the speech most closely perceived as a human voice.' (p. 10)	[35]
	'Naturalness refers to how closely the output sounds like human speech.' (p. 389.e1)	[42]
Combination of both	'By naturalness, we understand the voice stimulus to be perceived as a plausible outcome of the human speech production system' (p. 1)	[74]
	'[...] voices which sound like they could come from an actual human being (which should be rated as more natural) and voices that sound more fictitious, such as a cartoon character or a monster (which should be rated as less natural).' (p.429)	[57]

<sup>a</sup>Note: definitions are all original quotes from the respective references. The full compilation of extracted definitions can be found in the supplemental information online. Note that the mapping of definitions to the conceptualization of naturalness was carried out by us and not the authors of the original publications.

[20,63,64]. In addition, it concerns the experimental designs and measurements, especially rating scales that 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. Furthermore, the precise acoustic properties of voice material often remain elusive, bearing a risk for potential undetected confounds. Finally, only a few studies provide measurements of 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 that 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

#### 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 to 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 article offers a conceptual framework that can be tailored to any empirical design, such as 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, and any additional processing that is applied.
- Wherever possible, provide stimulus examples. This is important because readers may have a clear idea of how a male versus 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 such as **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 that could have unintended side effects on perceived naturalness [74,116].



impression formation according to social psychological models of intergroup 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 might argue that this is not problematic, because the different disciplines simply have different interests and readerships. However, some intriguing commonalities and systematic patterns only emerge when pooling evidence from all available angles. For example, across synthetic, pathological, and acoustically manipulated voices, converging evidence emerges for a strong effect of pitch variation on perceived naturalness [14,26,74]. Furthermore, although several studies failed to find an **uncanny valley** [75] effect for synthetic voices [11,76], a recent study suggests it might exist for pathological voices [77]. This lack of exchange between research fields not only has precluded relevant insights but also has likely impeded the visibility and impact of voice naturalness research as a whole.

#### 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 lacks 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.

#### Toward a concise framework for voice naturalness

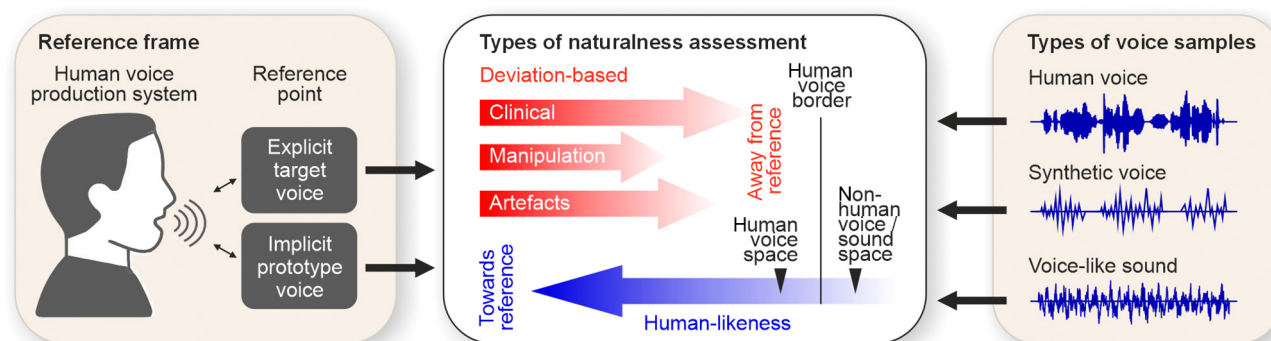
After identifying key problems that impede a systematic understanding of naturalness in voices, a logical next step is to propose concrete measures to address them, starting with a conceptual framework for the explicit definition of naturalness in voices.

#### Definitions of naturalness

We propose a taxonomy with two distinct types: **deviation-based naturalness** and human-likeness-based naturalness (Figure 2, Key figure). 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 for example, 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 for example, 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]).

**Key figure**

A conceptual framework for the definition of voice naturalness



Trends in Cognitive Sciences

**Figure 2.** 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 they lie within the perimeter of an acceptable human voice border.

**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 with 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 (Table 1), but they often remain implicit through certain design choices only (Box 1). For example, comparing human with 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 be not primarily a conceptual but an empirical matter requiring future inspection.

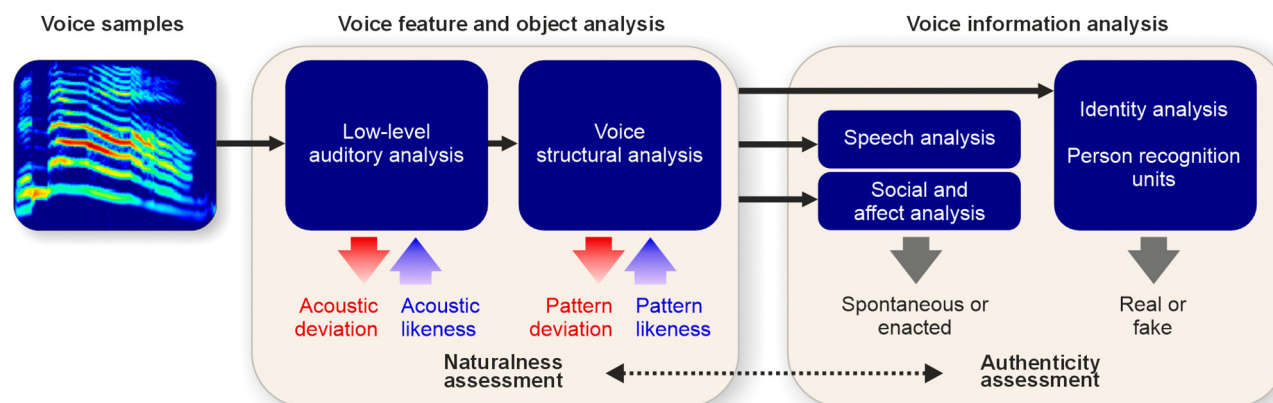
A second concept that deserves particular consideration is authenticity. In the scientific literature, authenticity is an established term with meaning that may refer to vocal emotion, identity, or gender – rather than the holistic impression of a voice. Emotional authenticity, for example, refers to 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. For example, ‘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 (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: (i) converting empirical heterogeneity from an impediment into an advantage and (ii) 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





Trends in Cognitive Sciences

**Figure 3. Rooting voice naturalness in voice-processing theory.** Theories of voice perception suggest a multilevel processing approach for voice samples (left panel), which involves analyzing these samples on the basis of 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 either can be expressed spontaneously (authentic) or can be enacted (nonauthentic), or it could be real or fake in relation to person-related identity information. Naturalness and authenticity assessments may have mutual influences.

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 not only will enhance mutual inspiration between clinicians and engineers but also could 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 benefiting quality of life [91]. For people who are predicted to lose their personal voice due to progressive disorders such as amyotrophic lateral sclerosis (ALS) or due to planned **laryngectomy**, current voice banking technology already allows 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 multilevel 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 and social content carried by the voice signal [4,93–95].

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

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

### Perspectives for future research

Our theoretical considerations on the processing of voice naturalness call for investigations of its time course and underlying brain mechanisms – relative to authenticity assessment but also to other voice characteristics. Initial evidence suggests that voice naturalness affects the brain response as early as 200 ms after voice onset and interacts with the processing of vocal emotions [99–101]. However, comparably early effects have been found for authenticity assessments [86,102,103]. Although the interpretability of these findings is limited by the potential influence of acoustic confounds, the findings suggest that naturalness and authenticity assessments are both fast and fundamental parts of voice perception. However, electrophysiological insights directly comparing the time course of naturalness and authenticity are elusive, as is their interplay with impressions of age, gender, or personality traits. A recent electroencephalographic (EEG) study suggests that many first impressions formed from voices are highly intercorrelated [8], but for naturalness, we are currently limited to behavioral data that point toward interactions with age, gender, and emotion perception [60,63,74]. In a broad sense, naturalness impressions are always formed against a specific context, whether that context refers to the voice itself or to the properties of the interaction. Accordingly, whether the same voice is assessed in an all-human or human–machine interaction context could make a crucial difference.

In that vein, while this review focuses on understanding naturalness in voices from an interdisciplinary perspective, we wish to emphasize the multisensory perspective of naturalness research. In fact, substantial research in the domain of faces has compared the perceived naturalness or realism of synthesized versus real faces (for a systematic review and meta-analysis, see [104]). Recent research even demonstrated conditions in which synthesized faces can be perceived as more human than genuine human faces. Moreover, an attempt to identify the visual features that trigger such a paradoxical facial ‘hyperrealism’ effect suggested contributions of typicality, familiarity, attractiveness, and low memorability [105]. Although this interpretation was based on qualitative reports and requires converging evidence, such research can inspire the systematic search for commonalities or differences between mechanisms that trigger judgments of voice or face naturalness. Ultimately, naturalness research should also systematically consider interactions between vocal and visual aspects of naturalness in combination. Indeed, accumulating

evidence suggests a complex interplay of visual appearance, vocal features, behavior, and the interactional context for the 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 that 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 the 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, composed 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 toward more systematic future research on numerous open questions regarding 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.

## Acknowledgments

We thank Simone Dahmen and Fatma Bilem for their support with the literature analysis and the members of the Jena Voice Research Unit (<https://www.voice.uni-jena.de/>) for helpful suggestions on this project. The authors gratefully acknowledge the award of funding through an EU's Marie Skłodowska-Curie Actions Doctoral Network ‘Voice Communication Sciences’ (action 101168998, <https://www.vocs.eu.com/>).

C.N.: I dedicate this work to our stillborn son. Thanks for changing our lives

## Declaration of interests

The authors declare no competing interests.

## Supplemental information

Supplemental information associated with this article can be found online at <https://doi.org/10.1016/j.tics.2025.01.010>.

## 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?

## References

- Román, S. *et al.* (2017) The importance of food naturalness for consumers: results of a systematic review. *Trends Food Sci. Technol.* 67, 44–57
- Meier, B.P. *et al.* (2019) Naturally better? A review of the natural-is-better bias. *Soc. Personal. Psychol.* 13, e12494
- Ode, A. *et al.* (2009) Indicators of perceived naturalness as drivers of landscape preference. *J. Environ. Manag.* 90, 375–383
- Young, A.W. *et al.* (2020) Face and voice perception: understanding commonalities and differences. *Trends Cogn. Sci.* 24, 398–410
- Rodero, E. and Lucas, I. (2023) Synthetic versus human voices in audiobooks: the human emotional intimacy effect. *New Media Soc.* 25, 1746–1764
- Rodero, E. (2017) Effectiveness, attention, and recall of human and artificial voices in an advertising story. Prosody influence and functions of voices. *Comput. Hum. Behav.* 77, 336–346
- Roswandowitz, C. *et al.* (2024) Cortical-striatal brain network distinguishes deepfake from real speaker identity. *Commun. Biol.* 7, 711
- Lavan, N. *et al.* (2024) The time course of person perception from voices in the brain. *Proc. Natl. Acad. Sci. U. S. A.* 121, e2318361121
- Lavan, N. (2023) How do we describe other people from voices and faces? *Cognition* 230, 105253
- Jiang, Z. *et al.* (2024) Comparison of face-based and voice-based first impressions in a Chinese sample. *Br. J. Psychol.* 115, 20–39
- Kühne, K. *et al.* (2020) The human takes it all: humanlike synthesized voices are perceived as less eerie and more likable. Evidence from a subjective ratings study. *Front. Neurobot.* 14, 593732
- Ilves, M. and Surakka, V. (2013) Subjective responses to synthesized speech with lexical emotional content: the effect of the naturalness of the synthetic voice. *Behav. Inform. Technol.* 32, 117–131
- Ilves, M. *et al.* (2011) The effects of emotionally worded synthesized speech on the ratings of emotions and voice quality. In *Affective Computing and Intelligent Interaction. ACII 2011. Lecture Notes in Computer Science* (6974) (D'Mello, S. *et al.*, eds), pp. 588–598, Springer
- Anand, S. and Stepp, C.E. (2015) Listener perception of monopitch, naturalness, and intelligibility for speakers with Parkinson's disease. *J. Speech Lang. Hear. Res.* 58, 1134–1144
- Moya-Galé, G. and Levy, E.S. (2019) Parkinson's disease-associated dysarthria: prevalence, impact and management strategies. *Res. Rev. Parkinsonism* 9, 9–16
- Damico, J.S., Ball, M.J., eds (2019) *The SAGE Encyclopedia of Human Communication Sciences and Disorders*, SAGE Publications
- Klopfenstein, M. *et al.* (2020) The study of speech naturalness in communication disorders: a systematic review of the literature. *Clin. Linguist. Phon.* 34, 327–338
- Frankford, S.A. *et al.* (2024) Contributions of speech timing and articulatory precision to listener perceptions of intelligibility and naturalness in Parkinson's disease. *J. Speech Lang. Hear. Res.* 67, 2951–2963
- Euler, H.A. *et al.* (2021) Speech restructuring group treatment for 6-to-9-year-old children who stutter: a therapeutic trial. *J. Commun. Disord.* 89, 106073
- Hardy, T.L.D. *et al.* (2020) Acoustic predictors of gender attribution, masculinity-femininity, and vocal naturalness ratings amongst transgender and cisgender speakers. *J. Voice* 34, 300.e11–300.e26
- Hyppa-Martin, J. *et al.* (2024) A large-scale comparison of two voice synthesis techniques on intelligibility, naturalness, preferences, and attitudes toward voices banked by individuals with amyotrophic lateral sclerosis. *Augment. Altern. Commun.* 40, 31–45
- Nass, C. *et al.* (1994) Computers are social actors. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems Celebrating Interdependence - CHI '94*, ACM Press
- Seaborn, K. *et al.* (2021) Voice in human-agent interaction. *ACM Comput. Surv.* 54, 1–43
- Triantafyllopoulos, A. *et al.* (2023) An overview of affective speech synthesis and conversion in the deep learning era. *Proc. IEEE* 1355–1381
- Schreibelmayr, S. and Mara, M. (2022) Robot voices in daily life: vocal human-likeness and application context as determinants of user acceptance. *Front. Psychol.* 13, 787499
- Baird, A. *et al.* (2018) The perception and analysis of the likeability and human likeness of synthesized speech. In *Interspeech 2018*, pp. 2863–2867, ISCA
- Lee, E.-J. (2010) The more humanlike, the better? How speech type and users' cognitive style affect social responses to computers. *Comput. Hum. Behav.* 26, 665–672
- Lu, L. *et al.* (2021) Leveraging 'human-likeness' of robotic service at restaurants. *Int. J. Hosp. Manag.* 94, 1–9
- Cambre, J. and Kulkarni, C. (2019) One voice fits all? *Proc. ACM Hum.-Comput. Interact.* 3, 1–19
- Eyssel, F. and Yanco, H. *et al.* (2012) If you sound like me, you must be more human. In *HRI' 12. Proceedings of the Seventh Annual ACM/IEEE Conference on Human-Robot Interaction: March 5-8, 2012 Boston, Massachusetts, USA*, pp. 125–126, Association for Computing Machinery
- Im, H. *et al.* (2023) Let voice assistants sound like a machine: voice and task type effects on perceived fluency, competence, and consumer attitude. *Comput. Hum. Behav.* 145, 107791
- McGinn, C. and Torre, I. (2019 - 2019) Can you tell the robot by the voice? An exploratory study on the role of voice in the perception of robots. In *2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pp. 211–221, IEEE
- Mitchell, W.J. *et al.* (2011) A mismatch in the human realism of face and voice produces an uncanny valley. *i-Perception* 2, 10–12
- Yorkston, K.M. *et al.* (1999) *Management of Motor Speech Disorders in Children and Adults*, PRO-ED
- Mawlim, C.O. *et al.* (2022) Speaker anonymization by modifying fundamental frequency and x-vector singular value. *Comput. Speech Lang.* 73, 1–17
- Hu, P. *et al.* (2021) Dual humanness and trust in conversational AI: a person-centered approach. *Comput. Hum. Behav.* 119, 106727
- Nusbaum, H.C. *et al.* (1997) Measuring the naturalness of synthetic speech. *Int. J. Speech Technol.* 2, 7–19
- Mayo, C. *et al.* (2011) Listeners' weighting of acoustic cues to synthetic speech naturalness: a multidimensional scaling analysis. *Speech Comm.* 53, 311–326
- Abdulrahman, A. and Richards, D. (2022) Is natural necessary? Human voice versus synthetic voice for intelligent virtual agents. *MTI* 6, 51
- Urakami, J. *et al.* (2020) The effect of naturalness of voice and empathic responses on enjoyment, attitudes and motivation for interacting with a voice user interface. In *Human-Computer Interaction. Multimodal and Natural Interaction* (Kurosu, M., ed.), pp. 244–259, Springer International Publishing
- Velner, E. *et al.* (2020) Intonation in robot speech. In *Proceedings of the 2020 ACM/IEEE International Conference on Human-Robot Interaction* (Belpaeme, T. *et al.*, eds), pp. 569–578, ACM
- Yamasaki, R. *et al.* (2017) Perturbation measurements on the degree of naturalness of synthesized vowels. *J. Voice* 31, 389.e1–389.e8
- Ko, S. *et al.* (2023) The effects of robot voices and appearances on users' emotion recognition and subjective perception. *Int. J. Human. Robot.* 20, 2350001
- Abur, D. *et al.* (2021) Feedback and feedforward auditory-motor processes for voice and articulation in Parkinson's disease. *J. Speech Lang. Hear. Res.* 64, 4682–4694
- Klopfenstein, M. (2015) Relationship between acoustic measures and speech naturalness ratings in Parkinson's disease: a within-speaker approach. *Clin. Linguist. Phon.* 29, 938–954
- Klopfenstein, M. (2016) Speech naturalness ratings and perceptual correlates of highly natural and unnatural speech in hypokinetic dysarthria secondary to Parkinson's disease. *JIRCD* 7, 123–146

47. Moya-Galé, G. *et al.* (2024) Perceptual consequences of online group speech treatment for individuals with Parkinson's disease: a pilot study case series. *Int. J. Speech Lang. Pathol.*, Published online May 1, 2024. <https://doi.org/10.1080/17549507.2024.2330538>
48. Eadie, T.L. and Doyle, P.C. (2002) Direct magnitude estimation and interval scaling of naturalness and severity in tracheoesophageal (TE) speakers. *J. Speech Lang. Hear. Res.* 45, 1088–1096
49. Eadie, T.L. *et al.* (2008) Influence of speaker gender on listener judgments of tracheoesophageal speech. *J. Voice* 22, 43–57
50. Yorkston, K.M. *et al.* (1990) The effect of rate control on the intelligibility and naturalness of dysarthric speech. *J. Speech Hear. Disord.* 55, 550–560
51. Schölderle, T. *et al.* (2023) Speech naturalness in the assessment of childhood dysarthria. *Am. J. Speech-Lang. Pathol.* 32, 1633–1643
52. Lehner, K. and Ziegler, W. (2022) Clinical measures of communication limitations in dysarthria assessed through crowdsourcing: specificity, sensitivity, and retest-reliability. *Clin. Linguist. Phon.* 36, 988–1009
53. Vogel, A.P. *et al.* (2019) Speech treatment improves dysarthria in multisystemic ataxia: a rater-blinded, controlled pilot-study in ARSACS. *J. Neurol.* 266, 1260–1266
54. Jones, H.N. *et al.* (2019) Auditory-perceptual speech features in children with Down syndrome. *Am. J. Intellect. Dev. Disabil.* 124, 324–338
55. Assmann, P.F. *et al.* (2006) Effects of frequency shifts on perceived naturalness and gender information in speech. In *INTERSPEECH*
56. Venkatraman, A. and Sivasankar, M.P. (2018) Continuous vocal fry simulated in laboratory subjects: a preliminary report on voice production and listener ratings. *Am. J. Speech-Lang. Pathol.* 27, 1539–1545
57. Kapolowicz, M.R. *et al.* (2022) Effects of spectral envelope and fundamental frequency shifts on the perception of foreign-accented speech. *Lang. Speech* 65, 418–443
58. Tamagawa, R. *et al.* (2011) The effects of synthesized voice accents on user perceptions of robots. *Int. J. Soc. Robot.* 3, 253–262
59. Mackey, L.S. *et al.* (1997) Effect of speech dialect on speech naturalness ratings: a systematic replication of Martin, Haroldson, and Triden (1984). *J. Speech Lang. Hear. Res.* 40, 349–360
60. Goy, H. *et al.* (2016) Effects of age on speech and voice quality ratings. *J. Acoust. Soc. Am.* 139, 1648
61. Coughlin-Woods, S. *et al.* (2005) Ratings of speech naturalness of children ages 8–16 years. *Percept. Motor Skills* 100, 295–304
62. Baird, A. *et al.* (2017) Perception of paralinguistic traits in synthesized voices. In *Proceedings of the 12th International Audio Mostly Conference on Augmented and Participatory Sound and Music Experiences* (Fazekas, G. *et al.*, eds), pp. 1–5, ACM
63. Merritt, B. and Bent, T. (2020) Perceptual evaluation of speech naturalness in speakers of varying gender identities. *J. Speech Lang. Hear. Res.* 63, 2054–2069
64. Baird, A. *et al.* (2018) The perception of vocal traits in synthesized voices: age, gender, and human likeness. *J. Audio Eng. Soc.* 66, 277–285
65. Aylett, M.P. *et al.* (2020) Speech synthesis for the generation of artificial personality. *IEEE Trans. Affect. Comput.* 11, 361–372
66. Kramer, R.S.S. *et al.* (2024) The psychometrics of rating facial attractiveness using different response scales. *Perception* 53, 645–660
67. Martin, R.R. *et al.* (1984) Stuttering and speech naturalness. *J. Speech Hear. Disord.* 49, 53–58
68. van Eck, N.J. and Waltman, L. (2010) Software survey: VOSviewer, a computer program for bibliometric mapping. *Scientometrics* 84, 523–538
69. van der Linden, S. (2023) *Foolproof: Why We Fall for Misinformation and How to Build Immunity*, W.W. Norton & Company
70. Fiske, S.T. (2018) Stereotype content: warmth and competence endure. *Curr. Dir. Psychol. Sci.* 27, 67–73
71. Todorov, A. *et al.* (2008) Understanding evaluation of faces on social dimensions. *Trends Cogn. Sci.* 12, 455–460
72. Sutherland, C.A.M. *et al.* (2013) Social inferences from faces: ambient images generate a three-dimensional model. *Cognition* 127, 105–118
73. Sutherland, C.A.M. *et al.* (2016) Integrating social and facial models of person perception: converging and diverging dimensions. *Cognition* 157, 257–267
74. Nussbaum, C. *et al.* (2023) Perceived naturalness of emotional voice morphs. *Cognit. Emot.* 1–17
75. Mori, M. *et al.* (2012) The uncanny valley. *IEEE Robot. Automat. Mag.* 19, 98–100
76. Romportl, J. (2014) Speech synthesis and uncanny valley. In *Text, Speech and Dialogue* (Horák, A. *et al.*, eds), pp. 595–602, Springer International Publishing
77. Diel, A. and Lewis, M. (2024) Deviation from typical organic voices best explains a vocal uncanny valley. *Comput. Hum. Behav. Rep.* 14, 100430
78. van Prooije, T. *et al.* (2024) Perceptual and acoustic analysis of speech in spinocerebellar ataxia type 1. *Cerebellum* 112–120
79. Moore, B.C.J. and Tan, C.-T. (2003) Perceived naturalness of spectrally distorted speech and music. *J. Acoust. Soc. Am.* 114, 408–419
80. Rao, M.V. *et al.* (2018) Effect of source filter interaction on isolated vowel-consonant-vowel perception. *J. Acoust. Soc. Am.* 144, EL95
81. Ratcliff, A. *et al.* (2002) Factors influencing ratings of speech naturalness in augmentative and alternative communication. *Augment. Altern. Commun.* 18, 11–19
82. Meltzner, G.S. and Hillman, R.E. (2005) Impact of aberrant acoustic properties on the perception of sound quality in electrolarynx speech. *J. Speech Lang. Hear. Res.* 48, 766–779
83. Andics, A. *et al.* (2010) Neural mechanisms for voice recognition. *Neuroimage* 52, 1528–1540
84. Valentine, T. *et al.* (2016) Face-space: a unifying concept in face recognition research. *Q. J. Exp. Psychol. (Hove)* 69, 1996–2019
85. Lima, C.F. *et al.* (2021) Authentic and posed emotional vocalizations trigger distinct facial responses. *Cortex* 141, 280–292
86. Sarzedas, J. *et al.* (2024) Blindness influences emotional authenticity perception in voices: behavioral and ERP evidence. *Cortex* 172, 254–270
87. Anikin, A. and Lima, C.F. (2017) Perceptual and acoustic differences between authentic and acted nonverbal emotional vocalizations. *Q. J. Exp. Psychol. (Hove)* 71, 622–641
88. Kachel, S. *et al.* (2020) Gender (conformity) matters: cross-dimensional and cross-modal associations in sexual orientation perception. *J. Lang. Soc. Psychol.* 39, 40–66
89. Mills, M. *et al.* (2017) Expanding the evidence: developments and innovations in clinical practice, training and competency within voice and communication therapy for trans and gender diverse people. *Int. J. Transgend.* 18, 328–342
90. von Efff, C.I. *et al.* (2022) Crossmodal benefits to vocal emotion perception in cochlear implant users. *iScience* 25, 105711
91. Schweinberger, S.R. and von Efff, C.I. (2022) Enhancing socio-emotional communication and quality of life in young cochlear implant recipients: perspectives from parameter-specific morphing and caricaturing. *Front. Neurosci.* 16, 956917
92. Yamagishi, J. *et al.* (2012) Speech synthesis technologies for individuals with vocal disabilities: voice banking and reconstruction. *Acoust. Sci. Technol.* 33, 1–5
93. Belin, P. *et al.* (2004) Thinking the voice: neural correlates of voice perception. *Trends Cogn. Sci.* 8, 129–135
94. Belin, P. *et al.* (2011) Understanding voice perception. *Br. J. Psychol.* 102, 711–725
95. Lavan, N. and McGettigan, C. (2023) A model for person perception from familiar and unfamiliar voices. *Commun. Psychol.* 1, 1–11
96. Staib, M. and Frühholz, S. (2023) Distinct functional levels of human voice processing in the auditory cortex. *Cereb. Cortex* 33, 1170–1185
97. Staib, M. and Frühholz, S. (2021) Cortical voice processing is grounded in elementary sound analyses for vocalization relevant sound patterns. *Prog. Neurobiol.* 200, 101982
98. Pinheiro, A.P. *et al.* (2021) Emotional authenticity modulates affective and social trait inferences from voices. *Philos. Trans. R. Soc. Lond. Ser. B Biol. Sci.* 376, 20200402



99. Duville, M.M. *et al.* (2022) Neuronal and behavioral affective perceptions of human and naturalness-reduced emotional prosodies. *Front. Comput. Neurosci.* 16, 1022787
100. Duville, M.M. *et al.* (2024) Improved emotion differentiation under reduced acoustic variability of speech in autism. *BMC Med.* 22, 121
101. Nussbaum, C. *et al.* (2022) Contributions of fundamental frequency and timbre to vocal emotion perception and their electrophysiological correlates. *Soc. Cogn. Affect. Neurosci.* 17, 1145–1154
102. Kosilo, M. *et al.* (2021) The neural basis of authenticity recognition in laughter and crying. *Sci. Rep.* 11, 23750
103. Conde, T. *et al.* (2022) The time course of emotional authenticity detection in nonverbal vocalizations. *Cortex; J. Dev. Study Nerv. Syst. Behav.* 151, 116–132
104. Miller, E.J. *et al.* (2023) How do people respond to computer-generated versus human faces? A systematic review and meta-analyses. *Comput. Hum. Behav. Rep.* 10, 100283
105. Miller, E.J. *et al.* (2023) AI hyperrealism: why AI faces are perceived as more real than human ones. *Psychol. Sci.* 34, 1390–1403
106. Cabral, J.P. *et al.* (2017) The influence of synthetic voice on the evaluation of a virtual character. In *Interspeech 2017*, pp. 229–233, ISCA
107. Ehret, J. *et al.* (2021) Do prosody and embodiment influence the perceived naturalness of conversational agents' speech? *ACM Trans. Appl. Percept.* 18, 1–15
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
110. Higgins, D. *et al.* (2022) Sympathy for the digital: influence of synthetic voice on affinity, social presence and empathy for photorealistic virtual humans. *Comput. Graph.* 104, 116–128
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. *Comput. Hum. Behav.* 141, 107645
112. Parmar, D. *et al.* (2022) Designing empathic virtual agents: manipulating animation, voice, rendering, and empathy to create persuasive agents. *Auton. Agent. Multi-Agent Syst.* 36, 1–24
113. Sarigul, B. and Urgen, B.A. (2023) Audio-visual predictive processing in the perception of humans and robots. *Int. J. Soc. Robot.* 15, 855–865
114. Lowry, H. *et al.* (2013) Behavioural responses of wildlife to urban environments. *Biol. Rev. Camb. Philos. Soc.* 88, 537–549
115. Kauk, J. *et al.* (2024) The adaptive community-response (ACR) method for collecting misinformation on social media. *J. Big Data* 11, 1–32
116. Malisz, Z. *et al.* (2020) *Modern speech synthesis for phonetic sciences: a discussion and an evaluation*. pp. 487–491