

# Vocal Tone Modulation in AI Voice Assistants: Enhancing User Trust and Engagement Through Adaptive Responses

Ulysse Duvillier – 910403  
u.duvillier@campus.unimib.it  
Università degli studi di Milano-Bicocca

## TABLE OF CONTENT

1. Introduction	2
2. The impact of latency in AI voice assistants	2
2.1. Textual assistant in contrast with vocal assistant	2
2.2. Psychological effects of latency	2
2.3. Voice modulation and the role of the Feedback Loops	2
3. The role of vocal tone in shaping perception of expertise	3
3.1. The connection between tone and expertise	3
3.2. Psychological effects of tonal modulation in multitasking scenarios	3
3.3. The importance of vocal attractiveness	4
4. Implementing adaptive vocal tone modulation in AI systems	4
4.1. Adaptive tone and clear speech strategies	4
4.2. Role-playing for perceived expertise	5
4.3. Unified latent space models for low latency	5
5. Broader implications and ethical considerations of vocal tone modulation	6
5.1. Applications in accessibility and inclusivity	6
5.2. Cross-cultural adaptability	6
5.3. Ethical risks and challenges	7
5.4. Future directions for research	7
6. Conclusion	7

## 1. INTRODUCTION

In the era of ubiquitous technology, artificial intelligence (AI) voice assistants such as Alexa, Siri, and Google Assistant have revolutionized Human-Computer Interaction (HCI). By offering users a natural and intuitive way to communicate with machines, these assistants have become an integral part of daily life, supporting tasks ranging from setting reminders to answering complex queries. Despite their growing popularity, these systems face a significant challenge: latency. Delays between user commands and AI responses can disrupt the flow of interaction, leading to user frustration and diminished trust in the technology.

To address this issue, researchers and developers have explored innovative approaches to managing latency without compromising user experience. One promising strategy lies in vocal tone modulation, an approach inspired by human communication. By dynamically adjusting tone and style, AI voice assistants can convey shifts in expertise or context, thereby enhancing user trust, managing expectations, and mitigating frustration during high-latency scenarios.

This essay explores how vocal tone modulation serves as an effective tool for AI voice assistants to simulate expertise, adapt to user expectations, and maintain engagement, ultimately improving interactions in the face of latency challenges.

## 2. THE IMPACT OF LATENCY IN AI VOICE ASSISTANTS

Latency, defined as the delay between user input and system response, is a critical challenge in the design of AI voice assistants. In Human-Computer Interaction (HCI), users expect near-instantaneous responses, particularly in conversational interfaces, where delays disrupt the natural flow of communication. Research shows that even slight delays can lead to frustration, reduced trust, and a perception of inefficiency in the system. For example, Park and Kim emphasize the importance of real-time feedback in maintaining user satisfaction when delays occur (Park and Kim).

Please consider that the latency is a key feature in maximizing the satisfaction of the user, which is the goal in ai assistant in general.

### 2.1. Textual assistant in contrast with vocal assistant

First, a key distinction between voice assistants and their textual counterparts lies in how users perceive and tolerate latency. With textual interfaces, users often receive immediate visual feedback, such as a loading indicator or partial response, signaling that the system is processing their request. This continuous feedback reassures the user and reduces frustration.

In contrast, vocal assistants lack this intermediate feedback during spoken exchanges, making latency periods feel like conversational “blanks.” Such silences can disrupt the user’s flow of interaction and create uncertainty about whether the system has understood or is responding to the input.

### 2.2. Psychological effects of latency

As we said, the psychological effects of latency are deeply rooted in user expectations. When interacting with AI systems, users assume their commands will be processed quickly and accurately. Delays can create a sense of uncertainty and diminish confidence in the assistant’s capabilities. The study by Lee, Cheon, and Wang explores how users’ psychological states, such as perceived trust and satisfaction, are influenced by system responsiveness and feedback mechanisms. They highlight that adaptive strategies like **vocal modulations** are effective in mitigating frustration caused by latency (Lee et al.).

### 2.3. Voice modulation and the role of the Feedback Loops

Timely feedback is crucial in mitigating these negative effects. Cohn and Zellou demonstrate that prosodic adjustments, such as **changing pitch or tone** to convey effort or focus, can effectively manage user expectations during latency. For instance, a calm, steady tone might indicate that the system is

“thinking,” thereby reducing the likelihood of frustration and improving overall user experience (Cohn and Zellou).

Additionally, Park and Kim underline that clear communication of system status through feedback loops fosters trust by signaling that the system is actively processing the user’s request (Park and Kim). The addition of exclamation such as “mmh” or “ooh let me think” adds a time delay in order to mask the generation time of the models. This give also more to the user a **human like feeling** to the conversation.

#### 2.3.1. *Works Cited*

- Cohn, R., and G. Zellou. “**Prosodic Adjustments in Voice-AI Interactions: Intelligibility and Speech Patterns.**” *Frontiers in Psychology*, 2021.
- Lee, J., Y. Cheon, and M. Wang. “**A Study of the Interaction Between User Psychology and Perceived Value of AI Voice Assistants from a Sustainability Perspective.**” *Sustainability*, vol. 15, no. 14, 2023, p. 11396. <https://doi.org/10.3390/su151411396>.
- Park, D., and E. Kim. “**Method of Interacting Between Humans and Conversational Voice Agent Systems.**” *Journal of Human-Computer Interaction*, 2023.

### 3. THE ROLE OF VOCAL TONE IN SHAPING PERCEPTION OF EXPERTISE

Humans instinctively adjust their vocal tone to convey authority, friendliness, or confidence during communication. Research shows that these tonal changes strongly influence how others perceive a speaker, shaping their impressions of trustworthiness or expertise. AI voice assistants can use this principle to build a sense of authority and reliability in user interactions.

#### 3.1. The connection between tone and expertise

Research by Krämer et al. highlights that variations in vocal tone can significantly affect how listeners judge authority and expertise. Lower-pitched voices, for instance, are commonly associated with competence and trustworthiness, particularly in professional or advisory contexts (Krämer et al.).

Similarly, von Trott and Krämer found that when voice assistants use **deeper vocal frequencies**, users tend to perceive them as more authoritative and capable. These tonal indications align user expectations with the AI system’s intended role, making the interaction more effective and reassuring.

Additionally, this connection is particularly important for vocal AI systems aiming to provide expert guidance, such as in healthcare or technical support. By using a deeper, confident tone in these contexts, the AI assistant reinforces its perceived credibility and improves user engagement.

Therefore, improving the perception of expertise can lead to catastrophic accident in real world scenario such as healthcare. **Training and disclaimers** should be introduce for staff using technical ai assistant in sensible fields.

#### 3.2. Psychological effects of tonal modulation in multitasking scenarios

When multitasking or experiencing delays, users often rely on tonal cues to interpret the assistant’s attentiveness and reliability. A calm and deliberate tone can signal that the assistant is “thinking” or carefully processing the user’s input, which helps reduce frustration (Krämer et al.). Zhao et al. describe how tonal adjustments, such as variations in pitch and pacing, can make voice **assistants seem more approachable and intelligent**.

These effects are amplified in high-stress scenarios, where users may be juggling multiple tasks. A soothing tone helps lower cognitive load and creates a sense of emotional stability, making the interaction feel more supportive.

Additionally, tonal shifts like slight pauses or deliberate intonation mimic natural human conversational habits, enhancing the assistant’s relatability and masking any underlying latency issues as we previously said (Cohn and Zellou).

### 3.3. The importance of vocal attractiveness

Another key aspect of tonal modulation is vocal attractiveness, which combines clarity, warmth, and modulation to create a pleasant listening experience. Krämer et al. found that users are more likely to **associate attractive voices with intelligence and competence**, making these characteristics essential for AI assistants aiming to provide an authoritative presence.

This aspect is particularly relevant when engaging users over long periods, as a monotonous or robotic voice can lead to disengagement. Dynamic tonal adjustments not only enhance user satisfaction but also make the assistant more memorable and engaging.

By leveraging tonal modulation to convey authority, attentiveness, and warmth, AI voice assistants can align user perceptions with their intended roles. These strategies enhance the user’s trust and engagement, particularly in contexts where reliability and expertise are essential.

#### 3.3.1. Works Cited

- Krämer, Nicole, et al. *“Voice of Authority: Professionals Lower Their Vocal Frequencies When Giving Expert Advice.”* Journal of Language and Social Psychology, vol. 38, no. 3, 2019, pp. 331-349. <https://doi.org/10.1007/s10919-019-00307-0>.
- von Trott, Nicolas, and Nicole Krämer. *“Human Behavior Research on Vocal Modulation: Impacts on Authority Perception.”* Psychonomic Bulletin & Review, vol. 30, 2023, pp. 22-34. <https://doi.org/10.3758/s13423-023-02333-y>.
- Zhao, Fei, et al. *“Building Trust Through Voice: How Vocal Tone Impacts User Perception of Attractiveness of Voice Assistants.”* arXiv Preprints, 2023. <https://arxiv.org/pdf/2409.18941>.
- Cohn, R., and G. Zellou. *“Prosodic Adjustments in Voice-AI Interactions: Intelligibility and Speech Patterns.”* Frontiers in Psychology, 2021. <https://www.readcube.com/articles/10.3389/fcomm.2021.675704>

## 4. IMPLEMENTING ADAPTIVE VOCAL TONE MODULATION IN AI SYSTEMS

The successful implementation of adaptive vocal tone modulation in AI voice assistants hinges on both advanced technical frameworks and a nuanced understanding of human communication. Neural synthesis models and real-time processing techniques enable dynamic tone adaptation, allowing AI systems to tailor their vocal outputs to specific contexts. These innovations bridge the gap between technical functionality and behavioral insights, fostering more engaging and intuitive user experiences.

### 4.1. Adaptive tone and clear speech strategies

The success of adaptive vocal tone modulation in AI systems relies on advanced technical frameworks and practical strategies that enhance user engagement and comprehension. Neural synthesis models, as highlighted by Zhao et al., form the backbone of real-time tone adaptation. These models dynamically adjust pitch, cadence, and timbre based on user input or environmental factors, allowing AI systems to signal cognitive effort and maintain user engagement even during high-latency scenarios (Zhao et al.).

Consider a customer service AI assistant designed to handle user inquiries about a delayed shipment. When the user expresses frustration, the neural synthesis model can detect emotional insights in the input (e.g., tone, word choice) and generate a comforting response with a calm and empathetic tone.

The assistant might say, “I understand how frustrating this must be. Let me check on this for you right away,” in a slower, steady cadence with softer pitch.

On the other hand, if the user’s input is neutral, such as asking for order tracking details, the assistant can adapt its tone to be clear and efficient without consideration on being extra careful and the voice pace of the answer can be speed up.

In addition to dynamic tone adaptation, clear speech strategies further enhance intelligibility and trust. Research by Cohn and Zellou demonstrates that prosodic adjustments, such as varying speech rate and pitch, improve listener comprehension, particularly in noisy or complex environments (Cohn and Zellou). These adjustments allow AI systems to signal attentiveness, making interactions feel more human and supportive, regardless of external challenges.

#### 4.2. Role-playing for perceived expertise

In line with these technical strategies, Jean-Louis Quéguiner, founder of Gladia, explores how role-playing and tone modulation can influence user perceptions of AI competence. He explains that AI voice assistants can simulate conversations between multiple agents, each representing a **distinct area of expertise**. By adjusting their tone according to the role of each agent, the AI can enhance the perceived depth and reliability of its responses. This dynamic vocal modulation helps users interpret the AI’s tone as a reflection of its cognitive focus, allowing the system to maintain trust even during multitasking or delays (Quéguiner).

For example in the video, Mr. Quéguiner emphasize on how a vocal AI agent respond to a specific query from the user on a company. During the oral chat, the AI assistant interpret two different persons with two different voices : one is an expert of countability with a deep voice and the other one, expert of law with an mid-tone voice.

#### 4.3. Unified latent space models for low latency

Modern AI systems rely on architectural innovations to achieve low latency while maintaining high-quality performance. Unified latent space models, like those used in OpenAI’s ChatGPT, represent a significant step forward in this domain. Unlike traditional architectures that handle tasks such as voice-to-text, text processing, and text-to-voice in discrete stages, unified models process everything within a single, interconnected space (Hai, Jiarui, et al. ). This integration minimizes the need for task-specific models, reducing processing delays and enhancing conversational fluidity. These voice, unified latent space, models are called **end-to-end speech models**.

One example that demonstrates the importance of low latency is a live performance scenario where a vocal AI assistant sings in a duo with a human. For such collaborations, precise timing is crucial to produce a harmonious result. Unified latent space models ensure that the AI can synchronize its responses to match the timing of the human performer, showcasing how these models excel in real-time interactions. Those example not be perform with a classical architecture such speech to text, process and then text to speech.

To further optimize performance and satisfaction, OpenAI employs techniques like “chain-of-thought reasoning” (Quéguiner). This method breaks complex queries into sequential steps, allowing the model to reason through each stage without retraining the system. For instance, the system can handle a multi-step query such as responding to a user’s technical request—by breaking it into manageable sub-problems, maintaining conversational coherence while minimizing latency.

In fact, the development of these systems is supported by datasets from platforms like youtube, as seen with OpenAI’s Whisper. Whisper was designed to generate high-quality audio training data by leveraging the rich content of youtube videos, which include metadata such as speaker age, gender, and regional

accents. Moreover, those vocal assistant are really performant as youtube is full of “teach-student like conversation of expert” (Quéguiner). We can give the example of podcasts that are really similar to a discussion with an AI vocal assistant. Different interlocutor with questions and precise answers with a vast variety of accent etc.

Whisper as never been only a gift from OpenAI to the open-source community but first a tool for their dataset. These datasets enable the creation of robust, contextually aware models capable of dynamic tone adaptation and precise speech synthesis.

However, while unified models excel in conversational and creative tasks, they may face limitations in structured enterprise environments where accessing and processing textual databases is necessary and where the “vocal latency space” can interage with more classical APIs. For such use cases, traditional task-specific systems might still be required to ensure accuracy and efficiency.

#### 4.3.1. *Works Cited*

- Cohn, R., and G. Zellou. “**Prosodic Adjustments in Voice-AI Interactions: Intelligibility and Speech Patterns.**” *Frontiers in Psychology*, 2021. <https://www.frontiersin.org/articles/10.3389/fpsyg.2023.1131222/full>
- Zhao, Fei, et al. “**Building Trust Through Voice: How Vocal Tone Impacts User Perception of Attractiveness of Voice Assistants.**” *arXiv Preprints*, 2023. <https://arxiv.org/pdf/2201.02792>.
- Quéguiner, Jean-Louis. **How Changing Vocal Tone and Role Play Influence the Perception of Expertise in AI Voice Assistants.** YouTube, 2022. [https://www.youtube.com/watch?v=KQrxq4QXhg&t=1225s&ab\\_channel=Underscore\\_](https://www.youtube.com/watch?v=KQrxq4QXhg&t=1225s&ab_channel=Underscore_)
- Hai, Jiarui, et al. “**EzAudio: Enhancing Text-to-Audio Generation with Efficient Diffusion Transformer**”, *arXiv Preprints*, 2024. <https://arxiv.org/html/2409.10819v1>.

## 5. **BROADER IMPLICATIONS AND ETHICAL CONSIDERATIONS OF VOCAL TONE MODULATION**

Vocal tone modulation in AI systems extends beyond improving latency and user satisfaction, with applications that raise exciting possibilities as well as ethical challenges. By tailoring vocal outputs to user preferences and contextual demands, these systems can enhance personalization, engagement, and accessibility in a wide range of domains. However, the increasing sophistication of such technologies also necessitates careful consideration of their societal impact.

### 5.1. **Applications in accessibility and inclusivity**

One transformative area for vocal tone modulation is improving accessibility. Prosodic adjustments, such as varying speech rate or pitch, can significantly enhance intelligibility for individuals with auditory processing difficulties. Cohn and Zellou demonstrate how these adjustments make interactions more accessible, especially in environments with background noise or for users who benefit from explicit vocal indications (Cohn and Zellou). Similarly, adaptive tones that convey empathy could be particularly supportive such as trisomic individuals, helping them feel understood and engaged during interactions.

### 5.2. **Cross-cultural adaptability**

As AI voice assistants are adopted globally, the ability to adapt tone for different cultural norms is

increasingly important. In some cultures, formal and authoritative tones are more effective, while others may favor warmth and friendliness. Research by Zhao et al. highlights that culturally aligned tonal adjustments can enhance trust and acceptance in diverse populations, making systems more inclusive (Zhao et al.). Integrating vocal tone modulation with multilingual capabilities offers the potential for AI assistants to better cater to these cultural preferences, ensuring relevance and effectiveness across different regions.

### 5.3. Ethical risks and challenges

While the benefits of vocal tone modulation are evident, the technology also poses ethical risks. Overly persuasive or emotionally manipulative tones could exploit user vulnerabilities, especially in sensitive areas like mental health or online shopping. For instance, systems using confident tones in healthcare scenarios may unintentionally overstate the accuracy of their advice, leading to over-reliance or misuse of AI recommendations (Krämer et al.). Moreover, the increasingly human-like nature of AI voices could blur boundaries between human and machine, raising concerns about deception and transparency.

The diversity within the dataset used to train AI models presents a significant challenge in ensuring fair representation and the inclusion of minority groups. Imbalanced data can lead to underrepresentation or misrepresentation of certain ethnic and minority groups, resulting in biased outputs that may harm individuals or communities (Lee et al.). Addressing these disparities is crucial to creating equitable AI systems that respect and uphold the dignity of all users.

### 5.4. Future directions for research

To address these challenges, future research should focus on ensuring transparency in AI interactions. Explicit disclaimers during conversations can help users understand the system’s limitations and capabilities, preventing over-reliance. Zhao et al. suggest combining vocal modulation with visual or textual cues to create balanced, multimodal user interactions that are engaging but not manipulative. Additionally, incorporating ethical design principles into AI voice systems could help mitigate risks while maximizing their potential to serve diverse user needs.

#### 5.4.1. Works Cited

- Cohn, R., and G. Zellou. **“Prosodic Adjustments in Voice-AI Interactions: Intelligibility and Speech Patterns.”** *Frontiers in Psychology*, 2021. <https://www.frontiersin.org/articles/10.3389/fpsyg.2023.1131222/full>
- Krämer, Nicole, et al. **“Voice of Authority: Professionals Lower Their Vocal Frequencies When Giving Expert Advice.”** *Journal of Language and Social Psychology*, vol. 38, no. 3, 2019, pp. 331–349. <https://doi.org/10.1007/s10919-019-00307-0>.
- Zhao, Fei, et al. **“Building Trust Through Voice: How Vocal Tone Impacts User Perception of Attractiveness of Voice Assistants.”** *arXiv Preprints*, 2023. <https://arxiv.org/pdf/2409.18941>.
- Lee, J., Y. Cheon, and M. Wang. **“A Study of the Interaction Between User Psychology and Perceived Value of AI Voice Assistants from a Sustainability Perspective.”** *Sustainability*, vol. 15, no. 14, 2023, p. 11396. <https://doi.org/10.3390/su151411396>.

## 6. CONCLUSION

In the dynamic field of Human-Computer Interaction, managing latency in AI voice assistants is a pressing challenge with profound implications for user satisfaction and trust. This essay has demonstrated how vocal tone modulation offers a powerful solution, enabling AI systems to signal expertise, reliability,

and empathy. By dynamically adapting tonal characteristics, these systems mitigate frustration, enhance engagement, and maintain seamless interactions during high-latency scenarios.

The research highlights the psychological and practical value of tonal modulation. Studies show that lower frequencies can establish authority and trustworthiness, while adaptive strategies such as prosodic adjustments help maintain clarity and user attention in challenging environments (Krämer et al.; Zhao et al.; Cohn and Zellou). Additionally, innovations like OpenAI's Whisper and the application of chain-of-thought reasoning illustrate the technical advancements driving these capabilities. These systems rely on diverse datasets and unified architectures, bridging the gap between complex computations and human way of thinking design.

Looking ahead, the development of AI voice assistants must balance technological innovation with ethical responsibility. While inclusivity and cross-cultural adaptability offer significant benefits, challenges like biased training data and the over use of persuasive tones require careful examination. Transparent and user-focused AI design will be crucial to ensure these systems remain equitable, empathetic, and effective in diverse contexts.

By integrating the principles of human communication with cutting-edge technological frameworks, AI voice assistants can evolve into trusted and intuitive companions. Prioritizing ethical considerations and fostering user trust will not only enhance user satisfaction but also ensure these systems contribute positively to a rapidly advancing digital world.

---

The number of words in this document is 3225 and there are 19671 letters.