# Anonymization of Voices in Spaces for Civic Dialogue: Measuring Impact on Empathy, Trust, and Feeling Heard

WONJUNE KANG*, Massachusetts Institute of Technology, USA

MARGARET A. HUGHES*, Massachusetts Institute of Technology, USA

DEB ROY†, Massachusetts Institute of Technology, USA

Anonymity is a powerful component of many participatory media platforms that can afford people greater freedom of expression and protection from external coercion and interference. However, it can be difficult to effectively implement on platforms that leverage spoken language due to distinct biomarkers present in the human voice. In this work, we explore the use of voice anonymization methods within the context of a technology-enhanced civic dialogue network based in the United States, whose purpose is to increase feelings of agency and being heard within civic processes. Specifically, we investigate the use of two different speech transformation and synthesis methods for anonymization: voice conversion (VC) and text-to-speech (TTS). Through a series of two studies, we examine the impact that each method has on 1) the empathy and trust that listeners feel towards a person sharing a personal story, and 2) a speaker's own perception of being heard, finding that voice conversion is an especially suitable method for our purposes. Our findings open up interesting potential research directions related to anonymous spoken discourse, as well as additional ways of engaging with voice-based civic technologies.

CCS Concepts: • **Human-centered computing** → **Interaction paradigms**; **Systems and tools for interaction design**.

Additional Key Words and Phrases: Voice Anonymization; Civic Dialogue Network; Participatory Media; Voice Conversion; Text-to-Speech Synthesis; Speech Interfaces

## 1 Introduction

In recent years, civic engagement has been greatly driven by the use of digital and Internet-based participatory media. Such participatory media include platforms such as social network services, blogs, podcasts, digital storytelling, and virtual communities. Although all of these media have distinct characteristics, according to Rheingold [61], they share the following commonalities. First, they enable all participants to broadcast as well as receive information in many different modalities, including text, images, audio, and video. Second, their value and power derives from the active

---

*Both authors contributed equally to this research.

†Deb Roy is part-time unpaid CEO of Cortico.

Authors' Contact Information: Wonjune Kang, wjkang@mit.edu, Massachusetts Institute of Technology, Cambridge, MA, USA; Margaret A. Hughes, mhughes4@mit.edu, Massachusetts Institute of Technology, Cambridge, MA, USA; Deb Roy, dkroy@mit.edu, Massachusetts Institute of Technology, Cambridge, MA, USA.

participation of many people, and the links between users form a public as well as a market. Third, when amplified by information and communication networks, they enable broader, faster, and lower cost coordination of a wide range of activities. The key power to this form of media lies in its participatory potential, as it affords people the opportunity to express their own voice in such a way that allows them to turn their self-expression into a form of public participation [3]. In particular, the public voices of individuals, when aggregated and put into dialogue with the voices of others, make up a fundamental component of the "public sphere," as defined by the political philosopher Jürgen Habermas [32]. Given the power and freedom to influence policy, the public sphere can become an essential instrument of democratic self-governance [40].

In these settings, "voice" can be thought of as the unique style of personal expression that distinguishes one's communications from those of others [61]. At the same time, a person's literal voice has certain qualities that make it particularly effective for the purposes of participatory media described above. In addition to linguistic information, the human voice contains a significant amount of paralinguistic information such as emotion, emphasis, contrast, and focus [74]. It can also encode other elements of self-expression that may not be encoded by simple grammar or vocabulary, such as irony or sarcasm [34]. These factors make speech a much richer medium for communication than other modalities such as written text, and have contributed to rapidly growing interest in voice-based platforms such as Discord, as well as short-form video platforms such as TikTok, Instagram Reels, and YouTube Shorts, where voice can play a significant role in sharing information [29].

In many such participatory media platforms, providing *anonymity* to the voice can be a critical component of the overall system. Anonymity affords certain freedoms to participants in civic discourse by providing protection from coercion and interference [6]. Additionally, it can promote greater openness of discussion by allowing people to express and discuss potentially controversial opinions and arguments without risk to their image or career. One way of providing anonymity within social discourse is by following the so-called "Chatham House rule" [35], in which participants in a meeting are "free to use the information received, but neither the identity nor the affiliation of the speaker(s), nor that of any other participant, may be revealed." In practice, this is often done by simply taking a transcript or an abbreviated textual summary of the spoken content, censoring out any personally identifiable information, and further processing it for downstream usage of the information.

However, in settings where people are speaking and sharing powerful stories or life experiences, directly hearing their voices can be much more impactful than reading text summaries or transcripts; listening to someone tell a narrative can convey its essence much more effectively than simply reading a written version or transcript of the same story [46]. At the same time, speech contains a significant amount of personal information about the speaker, and it is a distinct biomarker by which a speaker's identity may be determined [51]. This poses the question: *How might we preserve the nuance, emotion, and other rich information encoded in speech while simultaneously providing anonymity to a speaker?*

Motivated by the above, this work explores the use of *voice anonymization* methods, where the objective is to preserve the qualities of an original speaker's voice that convey rich paralinguistic information while masking the speaker's identity. In particular, we study the application of voice anonymization in the context of a *civic dialogue network* based in the United States, which consists of collections of multi-party conversations in which participants engage in dialogue with one another and share stories and life experiences in relation to community or social issues. A key component of the civic dialogue network is the release of particularly insightful or powerful stories to the broader public and committed leadership figures (see Section 3). However, there exist cases

where speakers may fear retribution or the release of sensitive information, which motivates the usage of a voice anonymization system that can mask their identities.

We consider two different speech transformation/synthesis methods for voice anonymization: voice conversion (VC) and text-to-speech (TTS). In particular, we utilize state-of-the-art neural network-based methods [39, 42] that enable high-quality manipulation and synthesis of voices in a way that was not possible only a few years ago. We perform two studies to measure the impact of voice anonymization within our civic dialogue network. In the first study, we survey participants using a crowdsourcing website ($n = 1500$) to investigate how anonymization using VC and TTS affects the perceptions of listeners towards a storyteller, especially in terms of empathy and trust. We also explore whether informing listeners that audio has been altered for the purposes of anonymization affects their responses. In the second study, we investigate speakers' perceptions of their own stories that have been anonymized via VC or TTS by recruiting actual participants in the civic dialogue network ($n = 21$). We perform anonymization on audio clips of stories they shared during real-life dialogues and survey how anonymizing speech affects the speakers' impressions of how their emotions and intended message are conveyed. We also study the extent to which they feel their identity has successfully been masked, and whether anonymization changes the different parties in their communities that they would be comfortable sharing their story with.

While a large body of prior research has studied the qualities of synthetic voices in human-to-speech interface settings [15], to the best of our knowledge, there has been no work studying the impact of using these technologies to alter how real peoples' voices and stories are conveyed. Our work contributes to the literature by exploring not just a novel use of speech transformation/synthesis technologies, but also by studying how they might be used to foster and enable greater civic discourse. While we believe voice and anonymization in democratic participation to be relevant for exploration in democracies globally, due to the civic dialogue network we consider being positioned in the United States, our work primarily focuses on this specific context.

## 2   Context and Related Work

### 2.1   The Role of Vocal Qualities in Spoken Language

In spoken language, much of the meaning is determined by context—that is, the objects or entities which surround a focal communicative event [74]. This means that the truth or validity of a proposition is determined by commonsense reference to experience. Consequently, spoken language tends to convey subjective information, and an important aspect of it is the establishment of a relationship between the speaker and the audience [74]. This contrasts with written language, in which there is a greater emphasis on logical and coherent argument and most of the meaning is provided directly by the text itself. These properties of speech make it an ideal medium for people to share content such as stories and life experiences. Speech allows a listener to pick up on paralinguistic cues that convey the valence of emotional experiences [45, 67], and it has been shown to communicate human-like mental capacities related to thinking and feeling [69].

Key to the listening experience is the role of *prosody*; that is, how a given speaker uses patterns of intonation, stress, and rhythm to provide variations in their speech. Prosody plays an important role when we listen to a voice, as it is a major determinant of expressiveness, emotion, and naturalness, which are some of the most important traits in human communication [30, 62]. It also makes up a significant component in human interaction [55]; prosody features have been shown to help segment discourse and provide an acoustic signal for listeners to better understand speech [63, 66]. Thus, by contributing to the expressiveness of a speaker's message, prosody can draw attention to and aid with understanding.

In the past, prosody representation in artificially generated voices was poor, as traditional speech synthesis methods were not able to match real human voices in terms of basic quality and naturalness. However, newer speech synthesis methods based on deep neural networks are now capable of generating speech that rivals that of humans in terms of naturalness and prosodic variation [14, 41–43]. This has led to the development of many more practically usable speech interfaces, such as personal voice assistants like Amazon's Alexa, Apple's Siri, or Google Assistant.

The rapidly increasing popularity and prevalence of such technologies have motivated the Human-Computer Interaction (HCI) and Computer-Supported Cooperative Work (CSCW) communities to extensively study speech interfaces [15], especially the the role of prosody and vocal qualities on human perception of both real and synthetic voices. For example, previous research has studied the role of expressivity on the perception of storytelling speech [49] or how people utilize voice quality and prosody in the identification of emotions and attitudes [31]. On the speech synthesis side, prior work has evaluated different types of TTS voices for long-form content [12] or compared TTS and human voices in a narrative advertising setting, measuring how modifying prosody affected listeners in terms of effectiveness, attention, concentration, and recall [64]. Other work has studied social implications and research challenges in designing voices for smart devices, treating voice assistants as social agents within a human-robot interaction framework [13].

Most of the above work addressed the qualities of synthetic and/or human voices in the context of human-to-human or human-to-speech interfaces. However, to the best of our knowledge, there has been no previous work studying the the human-perceived impact of using speech transformation or synthesis technologies to change how real peoples' voices are conveyed to an audience. Prior work on voice anonymization [21, 33, 53, 56, 57] has primarily focused on developing improved technical systems, with evaluation limited to "hard" metrics such as success at fooling speaker identification systems or mean opinion score (MOS) listening tests for quality and intelligibility [77]. In this work, we consider a novel setting and evaluation method for anonymization, looking more at its perceptual impact on speakers and listeners using more human-driven metrics such as empathy, trust, and the conveyance of emotions and meaning.

## 2.2 Voice, Participation, and Anonymity in Democracy

"Voice" is often used as a metaphor for participation within civics and democracies. Extensive political literature emphasizes the importance of participation in democracy, and explores various means through which civic actors may participate [5, 25]. In particular, we look to deliberation and discourse within democracies as a key form of political participation. Habermas outlines a foundational model of discourse within democracies through the "public sphere" [32]. In the public sphere, a public, or evolving group of people, gathers to discuss and debate through conversation issues of concern to them to form an opinion. These gatherings are social, can be broken into smaller groups to discuss various issues, and are intended to be rational with reasoned arguments. The public sphere described here requires that all people within a society have access to participation, and that they should be independent from any governing body or coercion. Furthermore, the public sphere should have the ability to hold the governing state to account.

However, theorists identify that this is not the case in any democracy globally; rather, we see an emergence of "subaltern counter publics," or other versions of "public spheres" that are often smaller, created by those marginalized within a society, and centered around a shared identity in which that community can gather and discuss issues important to them free from fear of retribution [4, 24]. Anonymization within modern democracy is as old as the protected vote, and is seen as another way to protect minority voices and opinions from retaliation by an oppressive or dominating group [48, 75]. In the United States, anonymity is prevalent throughout democracy, be it through voting systems, campaign funding, or political protest [6]. Especially for marginalized groups, or

in cases where specific insights might directly challenge some powerful governing body and that body wishes to suppress that knowledge, anonymity can promise some protection. It can give both positive and negative freedoms in that it protects civic actors from coercion and interference, and as well as in that it creates space for identity development or adjustment [6]. Yet, anonymity also comes with certain challenges in that it can allow citizens to share harmful or hateful speech without consequence [73]. Additionally, while in some cases anonymity can protect and uplift marginalized minority voices from retaliation, direct democracies with anonymous participation can also easily enable tyranny of the majority as well, with no social consequence for a majority continuously pursuing their own interests at the expense of the minority [75].

The science, technology, and society community has long explored the impact of anonymity and identity masking in digital spaces. Many observed that technology and the masking of identity in online spaces invites participants to construct a version of their identity that might be different from that in the offline world, having a liberating and positive impact, but also a potentially harmful one [8, 10, 19, 78]. The early days of online social networks saw a lively debate (which continues today) around identity, verification and trust, as well as the value and unique conditions that enable masked identity online. While some platforms were founded on the principle of known communities in which verified identity was at the root [10], other platforms, such as Yik Yak and Reddit, were designed for various levels of anonymity, resulting in a wide range of social dynamics in these spaces. Reddit has been used as a respite for the lonely and can provide supportive spaces for mental health discussions, while it can also create spaces that support socially taboo behavior [17, 80]. Meanwhile, on Yik Yak, design features such as ephemerality and hyper-locality have yielded different behavioral outcomes [68].

Many tools sit at this digital democratic intersection and invite anonymous participation in decision making and government decisions. Examples of digital tools that have been designed to enable e-government and civic engagement anonymously include Pol.is [1], CommunityPulse [37], and CommunityCrit [44]. Such tools have highlighted the value of community voice in informing civic decision making, and the unique ability of digital tools in gathering large amounts of community data, making sense of that data, and communicating patterns to inform decisions. Furthermore, in many high stakes context such as work environments, worker voices have been elevated to management to give feedback anonymously without retribution [2, 20].

We contribute to this literature by exploring anonymization of the literal voice in audio of spoken stories. We look at listeners' perceptions of storytellers when their voices have been anonymized or not, as well as whether knowledge about whether the speech utterance has been anonymized affects their impressions. Finally, we evaluate not the behavioral outcome of people as anonymous actors, but rather the impact on their feelings of being heard by others when their voices have been anonymized.

## 3 Application Setting: Technology-Enhanced Civic Dialogue Network

The work in this paper was done in the context of the Fora conversation platform developed by the nonprofit Cortico,[1] a *technology-enhanced civic dialogue network* set in the United States. In this section, we briefly overview the structure and value of the civic dialogue network and the role of voice anonymization within it.

Throughout the United States and globally, there has been a decline in trust, increased polarization, and threats to democracy [22, 54]. Some members of the public feel generally unheard, expressing dissatisfaction with the current political system as well as their agency and ability to meaningfully take part in broader social decisions that affect their lives [36]. The conditions of traditional means

---

[1]https://cortico.ai/

of civic participation such as voting and town halls often invite reductive, non-nuanced expressions of opinions that can restrict participation and perpetuate injustices [36]. Consequently, there have been many explorations into how to more effectively share peoples' voices and allow them to participate in civics with nuance and greater agency, in pursuit of a more just civic system and democracy. These explorations encompass a range of alternative forms of decision making, from citizens' assemblies [23] to community organizing and movement building [27, 28, 82] to technology-enhanced civic participation [9, 37, 44].

Fora contributes to this field of explorations by attempting to increase feelings of agency and being heard within civic processes. It aims to do so via the following: 1) facilitated dialogue meant to increase understanding and connection between participants through nuanced story, question, and opinion sharing, 2) recording and systematically analyzing voices and themes from those dialogues to reveal patterns across and within communities around key issues, and 3) partnering with civic leaders committed to listening and acting upon these insights. The platform is run by a central non-profit organization which supports and trains civic leaders, community organizations, corporate spaces, schools, etc. to host facilitated dialogues within their communities around specific issues. Facilitated dialogue is an age-old method of gathering in which a figure with some authority, a facilitator, supports a gathered group through an experience-centered, structured conversation. Usually, there are norms that structure turn taking, types of contributions, the design of the space, and how to listen and engage with others. The facilitator often guides the conversation with the aid of a script or developed notes to elicit nuanced experiences, opinions, and questions from participants. These types of dialogues are a key method used in restorative justice spaces to heal harms, perform mediation across conflicts and differences, and deepen already existing relationships through Circle practice [7, 50, 52].

While the conversations themselves are of great value, a key component of Fora is that it enables content shared within the conversations to be heard elsewhere as well. By recording the conversations, analyzing them with participatory, qualitative methods, and sharing the emergent patterns through public-facing portals, specific highlighted stories and key themes expressed in the dialogues can be explored by a much larger public audience. Furthermore, many conversation campaigns are held in partnership with civic leadership who commit to listening, engaging, and responding to these voices. This means that dialogues can have a direct and visible path to practical impact and influence, potentially increasing participant feelings of being heard and having agency.

Central to the design of the network is the human voice. Participants verbally share their stories, which are often tied to life experiences; thus, any data point in the system is always tied to one or more audio recordings of participants' voices. In particular, a key aspect of the system is that, when a highlighted story is shared in a public portal, the audio of the participant's voice is always shared along with the text transcript (this is only done if the participant grants consent). This is an important and central design choice for the network because of the reasons outlined in Section 2.1. However, there are settings in which people may be reluctant to allow their voices to be made publicly available if they have shared sensitive content or fear retaliation against what they have said, even if their story is especially powerful or insightful. Furthermore, as outlined in Section 2.2, the general capability for anonymity can be a highly valuable tool in social and political spaces to allow people to participate fully and authentically.

For these reasons, we explore how we might integrate voice anonymization into the Fora civic dialogue network. Specifically, we are interested in studying how to best maintain the rich information within the human voice, but transform it in a way that successfully hides the speaker's identity. Given that feelings of being heard and hearing others are key goals in the network, it is essential to evaluate how anonymization methods affect those goals or not. In the following sections, we outline our methods and results as we explore these questions through two studies.

## 4 Methods

We treat voice anonymization as the task of suppressing personally identifiable attributes of a speech signal while preserving the key attributes related to semantic content. It is relatively simple to anonymize speech by modulating voice characteristics using signal processing techniques. For example, one may alter the frequency spectrum to change the perceived pitch of the voice or design acoustic filters to change the spectral characteristics of the speech; many speech anonymization methods popularly used in the media today (e.g., for anonymous interviews) fall under these categories. However, these methods often result in robotic, unnatural sounding transformed voices. Because they are not explicitly designed to retain the comprehensibility of the speech, they can also necessitate the use of subtitles or captions to allow listeners to understand what is being said. In this study, we were specifically interested in voice anonymization methods that did not have this drawback, as perceptual intelligibility and clarity for listeners are critical components for our civic dialogue network setting. Therefore, we resorted to two speech transformation/synthesis technologies that satisfy our basic surface level requirement of preserving intelligibility: voice conversion (VC) and text-to-speech (TTS).

Voice conversion is the task of transforming a voice to sound like another person without changing the linguistic content of the original utterance [72]. It belongs to the general field of speech synthesis, which includes TTS as well as the changing of other speech properties such as emotion and accents. In our setting, we further define it as the task of changing a speech utterance's timbre (i.e., a speaker's vocal tone) while leaving all other aspects of the utterance such as prosody, rhythm, and accent unchanged [58]. VC has been used as the key component of many previously proposed voice anonymization systems [21, 33, 53, 56, 57]. Meanwhile, TTS is the task of producing synthetic speech that corresponds to the verbalization of a text input. It is the core technology that drives the voices of many speech interfaces such as Amazon's Alexa, Apple's Siri, and Google Assistant.

One might naturally expect VC to better fulfill the desired characteristics of a voice anonymization system for our civic dialogue network, as it preserves both the linguistic and prosodic content of a speech utterance while only changing the perceived speaker identity. Meanwhile, TTS only preserves linguistic content, removing prosodic characteristics and other paralinguistic information in addition to changing the speaker identity. However, we sought to compare these two "levels" of anonymization and study the impact that their differences would bring.

To perform VC, we used LVC-VC XL [39], a recently proposed deep neural network-based voice conversion model. The model was trained on the VCTK Corpus [81], which consists of around 44 hours of English speech from 109 speakers (62 female, 47 male) with various accents (American, Australian, Canadian, English, Indian, etc.). The model is *zero-shot*, which means that it is able to convert any *source* speaker's voice to sound like any *target* speaker's voice; this was a necessary characteristic because the model needed to be usable for any arbitrary speaker. For simplicity, we limited our selection of candidate target speaker voices to the 109 speakers that the model had seen during training and randomly sampled speakers from the pool of candidates when performing anonymization. To control for gender, we always selected target speakers that were of the same gender as the original speaker; further details on speaker selection are described in Sections 4.1 and 4.2. Note that LVC-VC only changes the source speaker's timbre—all other aspects of the source utterance, including rhythm, intonation, and accent, stay the same.

For TTS, we first performed automatic speech recognition (ASR) using Whisper-base [59] to obtain the text transcript of the speech utterance.[2] Then, we manually corrected any errors and

---

[2]The Whisper ASR model was downloaded and run locally on the authors' computing hardware, rather than using OpenAI's API.

formatted the transcript so as to match the rhythm and tempo of the original speech utterance as closely as possible (for example, adding a comma inserts a brief pause between words). Finally, to synthesize speech, we used VITS [42], a state-of-the-art multi-speaker neural network model, which was also trained on the VCTK Corpus. Here, we limited our pool of voices to speakers with American accents, specifically General American English accents, which left us with 17 female voice candidates and 5 male voice candidates (22 total). As with VC, we always performed anonymization by selecting speakers that were of the same gender as the original speaker. Note that TTS candidate voices were not specifically selected to match the accents of the original speakers; while most of the original speech utterances in our studies were spoken in General American English accents (7/10 in Study 1 and 18/21 in Study 2), some were not. It is possible that a perceived difference in speaker demographics could have affected listeners' perceptions, as has been observed and documented in prior studies [18, 26, 65]. However, it was unfortunately infeasible to completely control for these factors because of the limited set of accents that the TTS model could synthesize speech in (for example, it is not able to synthesize speech in Spanish or other Hispanic accents). We leave addressing this as potential future work. Further details on speaker selection for each study are described in Sections 4.1 and 4.2.

We conducted two studies to measure the impact of voice anonymization when done using either VC or TTS. In the first study ($n = 1500$), we investigated the impact of anonymization on *listeners* of stories told in our civic dialogue network, measured in terms of the empathy and trust they felt towards the storyteller. In the second study ($n = 21$), we investigated *speakers'* perceptions of their own stories in the civic dialogue network that had been anonymized, as well how effectively they felt that their identity was masked. Both study protocols were reviewed and approved by a university's institutional review board (IRB). In the rest of this section, we describe the design, procedures, and analyses done for each of the two studies.

## 4.1 Study 1: Perceptions of the Listener

In our first study, we investigated how anonymizing stories told in our civic dialogue network affected the perceptions of listeners towards the storyteller. We also considered how knowledge of whether a story had been anonymized or not (i.e., telling a listener that the speech utterance had been altered for the purposes of anonymization) affected their responses. Concretely, we aimed to answer the following research questions (RQs):

- **RQ 1:** How does anonymizing speech impact a listener's perception of and/or connection with the storyteller, especially in terms of trust and empathy?
- **RQ 2:** How does knowledge of whether speech has been altered for the purposes of anonymization affect a listener's trust and empathy towards the storyteller?

*4.1.1 Procedure.* We collected audio clips of 10 stories from our civic dialogue network that had been spoken by a variety of American speakers across age, gender, and emotional valence.[3] The stories covered topics such as housing stability, police brutality, trans rights, sibling relationships, and neighborhood violence. Five stories were spoken by women, while five were spoken by men. The stories ranged from 51 seconds to 2 minutes and 56 seconds in length ($\mu = 89.6$ seconds, $\sigma = 36.3$ seconds). More detailed information on each of the stories is shown in Table 1.

We generated anonymized versions of the stories by performing VC and TTS on the audio files as described above. For VC, we randomly selected ten potential target speakers of the same gender as the storyteller from the VCTK Corpus and used them to generate voice converted candidates.

---

[3]The exact age breakdown of the speakers was unavailable because we took the stories from conversations where that information was not collected.

Table 1. Topics, speaker genders, and durations of the original, VC, and TTS audio files for the 10 stories used in Study 1.

| Story # | Topic | Gender | Duration (Orig./VC) | Duration (TTS) |
|---|---|---|---|---|
| 1 | Rent, housing stability, and homelessness | Female | 2:10 | 1:06 |
| 2 | Police brutality on a family member | Female | 2:56 | 2:01 |
| 3 | Dental care for children | Female | 1:11 | 0:46 |
| 4 | Laws and support for trans rights | Female | 1:36 | 1:00 |
| 5 | Life lessons learned from a mother | Female | 1:17 | 0:47 |
| 6 | Parental pressure on children playing sports | Male | 0:58 | 0:47 |
| 7 | Growing up, maturing, and building relationships with siblings | Male | 1:39 | 1:06 |
| 8 | Neighborhood gun violence and its impact on youth | Male | 0:51 | 0:33 |
| 9 | Recognition of social privilege | Male | 1:01 | 0:58 |
| 10 | Volunteering for people from a different social background | Male | 2:10 | 1:06 |

Then, we manually selected one VC candidate that sounded sufficiently different from the original speaker's voice while maintaining clear audio quality. For TTS, we generated candidates using all of the American-accented voices from the VCTK Corpus that were of the same gender as the storyteller (17 for female, 5 for male) and manually selected the candidate that sounded the most natural. Audio that had been anonymized using VC maintained the same length as the original speech utterances. For TTS, the synthesized speech ranged from 33 seconds to 2 minutes and 1 second in length ($\mu$ = 61.2 seconds, $\sigma$ = 22.6 seconds); the shorter lengths reflected the lack of pauses, repetitions, and filler words that are part of regular human speaking patterns [11] but were not included in the generated TTS audio.

The objective of this study was to measure participants' responses towards the VC and TTS voices and compare them against responses towards the original, un-anonymized voice. Each study participant listened to a single random audio clip corresponding to one of three types of audio (Original, VC, or TTS) for one of the 10 selected stories. Participants were also randomly assigned to an "aware" or "unaware of anonymization" condition; those in the "aware" conditions saw the following message before being shown the audio player:

> Please note that the voice you will be hearing has been altered in order to anonymize the identity of the speaker.

while those in the "unaware" condition were only shown the audio player with no other text beyond the instructions. We included a check in the survey to ensure that all participants finished playing the audio before they could proceed. Then, participants answered a series of 20 Likert scale questions that were designed to measure various aspects of their trust and empathy towards the storyteller in the audio clip, as well as the following free response question:

> What were some aspects of the audio clip that made you provide the answers that you did above?

Further details on the survey design are described in Section 4.1.2.

Altogether, there were 6 audio type and knowledge conditions: 3 audio (Original, VC, and TTS), and 2 knowledge (aware, unaware).[4] For each condition on each of the 10 stories, we collected 25 responses. In total, this made for 1500 survey participants.

---

[4]Note that one combination of conditions, Original + aware, was a deceptive condition in which we told listeners that the audio had been altered although it was not.

Table 2. Outline of the 20 Likert scale survey questions used in Study 1. Each set of questions was developed to address a different relevant theme, namely empathy or trust.

| Theme | Question |
|---|---|
| Empathy | **E1:** I believe the speaker's emotions in the story were genuine.<br>**E2:** I experienced the same emotions as the speaker when listening to the story.<br>**E3:** I was in a similar emotional state as the speaker when listening to the story.<br>**E4:** I could feel the speaker's emotions.<br>**E5:** I could see the speaker's point of view.<br>**E6:** I recognized the speaker's situation.<br>**E7:** I could understand what the speaker was going through in the story.<br>**E8:** The speaker's reactions to the situation were understandable.<br>**E9:** When listening to the story, I was fully absorbed.<br>**E10:** I could relate to what the speaker was going through in the story.<br>**E11:** I could identify with the situation described in the story.<br>**E12:** I could identify with the characters in the story. |
| Trust | **T1:** I believe that the speaker was telling an authentic life experience.<br>**T2:** I believe the story is true.<br>**T3:** I felt like I could trust the speaker.<br>**T4:** I have respect for the speaker.<br>**T5:** I believe that the speaker wanted me to understand their experience and perspective.<br>**T6:** I believe the speaker was manipulating me or trying to persuade me with their story.<br>**T7:** I would like to learn more about the speaker.<br>**T8:** I would be willing to talk with the speaker and share some of my own experiences. |

*4.1.2 Survey design.* Given the key goals of our civic dialogue network, we looked to measure the empathy and trust of listeners towards storytellers after listening to audio of their stories. To measure this, we created a 5 point Likert scale survey (1 – strongly disagree, 5 – strongly agree) measuring *state empathy* and *trust* of the listener towards the storyteller. For empathy, we collected 12 survey questions from established research around state empathy [71]. We chose state empathy over other empathy measures as we were not interested in evaluating the general empathy of the listener, but rather their experience of empathy towards the speaker informed solely by the story they heard. For trust, we developed a set of 8 questions jointly with a key conversation designer within the civic dialogue network. We then validated with the conversation designer that both sets of questions on empathy and trust were measuring factors that were in alignment with the civic dialogue network's purposes and goals. The survey questions are shown in Table 2.

*4.1.3 Participants.* We recruited all participants on the Prolific crowdsourcing platform. Participants were required to be located in the United States or United Kingdom and be fluent in English. 53.33% of participants identified as female, 46.40% identified as male, and the remaining 0.27% preferred not to say. The racial breakdown of participants was 84.93% White, 5.8% Asian, 4.07% Black, 3.8% Mixed, and 1.4% Other. Participants were paid $1 USD for their time. The median time taken for the task was 4 minutes and 6 seconds, which resulted in a reward per hour of $14.63 USD.

## 4.2 Study 2: Perceptions of the Speaker

In our second study, we looked to measure speakers' perceptions of their own speech that had been anonymized via VC or TTS. The research questions were:

- **RQ 1:** How does anonymizing speech affect a speaker's own impression of how their emotions and intended message are conveyed?

- **RQ 2:** How do different speech anonymization methods affect a speaker's perception of their privacy?

*4.2.1    Participants and procedure.* To understand the impact of voice anonymization on speakers' perceptions of their own stories and the degree to which they felt heard, we reached out to a pool of participants from one set of conversations that had been held using the civic dialogue network described in Section 3. Given that the application of our research deals with the challenges of feeling agency and being heard within civic spaces, it was important that we looked to participants of real-world conversations that had actual consequences; without the weight of dealing with real stories, people may have responded superficially to our study. Furthermore, as a key aspect of voice anonymization is the protection of privacy, it was important that we choose participants and stories for which privacy was relevant to some degree. Therefore, we worked with a set of conversations that was held within a university's work environment, in which participants discussed both the joys and the frictions that were part of their workplace.

We reached out to all participants who took part in at least one of two such conversation campaigns within their workplace, which made up a pool of 82 people. A total of 21 participants completed our survey (see Section 4.2.2). 71.43% of participants identified as female, 23.81% identified as male, and 4.76% identified as non-binary. The racial breakdown of participants was 61.9% White, 19.05% Hispanic or Latino, and 9.52% Asian. Participants were compensated for their involvement in the study by being entered into a lottery with the chance to win a $50, $35, or $20 gift card.

*4.2.2    Survey design.* The survey for this study aimed to measure how storytellers perceived their own stories and voices that had been anonymized using VC and TTS. Specifically, we evaluated how anonymizing the audio of stories using the two methods affected storytellers' perceptions of conveyed emotions, their feelings of being heard, and their feelings of security and privacy when sharing their stories with others.

To do this, we first had participants listen to their own story in the original, unaltered voice. Then, they were asked to imagine a hypothetical scenario in which their stories would be shared through a public-facing digital portal that would be visible to all members of their broader workplace community, including their peers and superiors. Under this condition, they were told to consider the situation of sharing an anonymized version of their story on the portal. Participants then listened to both VC and TTS versions of their story that we generated using the following procedure. For VC, we generated four candidates per story using 2 female and 2 male voices randomly chosen from the VCTK Corpus as target speakers; we presented all four options to the participant and asked them to choose which version they would most like to use in the hypothetical scenario. For TTS, we followed the same procedure as in Study 1, where we generated candidates using the voices in American accents from the VCTK Corpus that were of the same gender as the participant. Then, we manually selected the single option whose speech sounded the most natural.

After listening to the VC and TTS candidates, participants answered the following 5 point Likert scale questions (1 – strongly disagree, 5 – strongly agree) in response to each type of audio:

(1) *I feel like my story would still be heard with the* (transformed, synthesized) *voice.*
(2) *I feel like the* (transformed, synthesized) *voice retained the authenticity of my story.*
(3) *I feel like the* (transformed, synthesized) *voice conveys the emotion I wanted to convey.*
(4) *I feel like the* (transformed, synthesized) *voice retains important characteristics of my voice (accent, speech pattern, intonation, etc.).*
(5) *I feel like the* (transformed, synthesized) *voice masked my identity.*

as well as the following free response question:

> *Please provide any additional comments on your impressions of the* (voice conversion, text-to-speech) *technology applied on your voice in the context of* [your community].

This was then followed by a question asking the participant about the members of their community with whom they would be comfortable sharing their story in the original, VC, and TTS voices. The question was posed as a checkbox selection question where multiple selections were possible, with the following options:

- **Group 1:** My community outside of my workplace (friends, acquaintances, family)
- **Group 2:** My peers (other graduate students, other staff members, whatever community you are a part of)
- **Group 3:** The broader university community
- **Group 4:** My supervisor(s)
- **Group 5:** University administration
- **Group 6:** People who have power over me
- **None of the above**

The survey finished with a question asking if their answers changed across the three conditions (Original, VC, and TTS), and if so, to elaborate on why.

### 4.3  Data Analysis Methods

*4.3.1  Likert scale questions.* We performed two-sided Student's $t$-tests to compare the responses to the Likert scale questions between the various conditions. For Study 1, we performed $t$-tests to compare responses for Original vs. VC and Original vs. TTS for all 20 questions. We also performed tests to compare the responses for each audio type (Original, VC, TTS) when listeners were aware vs. unaware of anonymization. For Study 2, we performed $t$-tests comparing the speakers' responses for the VC and TTS audio.

*4.3.2  Free response questions.* Because of the large number of responses to the free response question in Study 1, we performed a topic analysis of the answers rather than manually going through and analyzing them. To do this, we leveraged Sentence-BERT [60], a neural network language model that has been trained to encode sentences into semantically meaningful representations. Sentence-BERT maps sequences of words (usually a sentence or phrase) to 768-dimensional vector representations (henceforth referred to as "embeddings"). Intuitively, the model maps semantically similar sentences or phrases to embeddings that are close to each other and semantically different ones to embeddings that are farther from each other in the latent vector space.

We split all 1500 free response answers into sentences using the Python Natural Language Toolkit (NLTK) package, producing 2024 sentences. The full responses had an average length of 108 characters, and individual sentences had an average length of 80 characters. Then, we computed embeddings for the sentences and clustered them using the Python implementation of HDBSCAN[5] [47] with a minimum cluster size of 2. Setting the minimum cluster size to 2 ensured that any response that was not in a cluster on its own would be represented as a topic/theme. This resulted in 109 initial clusters formed by 928 responses (HDBSCAN also produced a null cluster containing 1096 elements). Then, we went through the clusters and agglomeratively merged similar ones until there remained 20 groups of responses that represented a distinct set of emergent topics and themes. We describe the results of this analysis in Section 5.1.3.

---

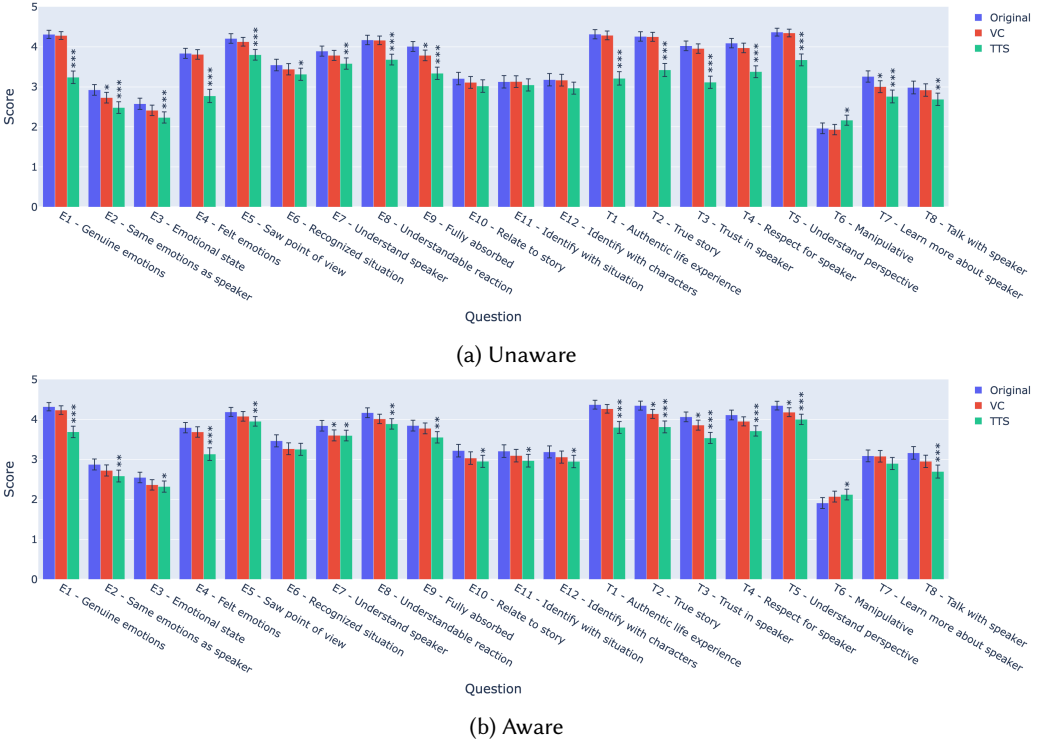[5]https://hdbscan.readthedocs.io/en/latest/

Fig. 1. Average scores for the 20 Likert scale questions in the survey for Study 1 when listeners were (a) unaware and (b) aware of anonymization. Error bars represent 95% confidence intervals. *, **, or *** on top of bars for VC and TTS denote statistically significant differences at $p < 0.05$, $p < 0.01$, or $p < 0.001$, respectively, for two-sided $t$-tests comparing against scores for Original.

## 5 Results

In this section, we summarize our findings from the two studies. Section 5.1 covers the results from Study 1, while Section 5.2 covers the results from Study 2. In each subsection, we first analyze quantitative results taken from the Likert scale questions, followed by a more qualitative analysis of the answers to the free response questions.

### 5.1 Study 1

*5.1.1 Listeners' perspectives on voice conversion vs. TTS for anonymization.* We first analyze the perceptions that listeners had towards stories that were anonymized using VC or TTS when compared against a baseline of the original, unaltered voice. Fig. 1 shows bar charts of the average scores for each of the 20 Likert scale questions in the survey, along with 95% confidence intervals, when listeners were (a) unaware and (b) aware of anonymization. Statistically significant differences between scores for the Original vs. VC/TTS conditions are denoted by asterisks above the bars. For figure compactness, we label questions on the $x$-axis with abbreviated codes (E1, E2, etc.); the mapping back to the full text of the question can be found in Table 2.

We find that stories anonymized using VC largely achieve similar scores to Original audio under both unaware and aware conditions. When listeners were unaware of anonymization, there were statistically significant differences ($p < 0.05$) between Original and VC audio for only 3/20 questions;

Fig. 2. Average scores for the 20 Likert scale questions in the survey for Study 1 comparing the unaware (blue) and aware (red) conditions for (a) Original, (b) VC, and (c) TTS voices. Error bars represent 95% confidence intervals. *, **, or *** on top of bars for Aware denote statistically significant differences at $p < 0.05$, $p < 0.01$, or $p < 0.001$, respectively, for two-sided $t$-tests comparing against scores for Unaware.

when listeners were aware of anonymization, there were statistically significant differences for only 4/20 questions. This indicates that VC is successful at meeting our objectives for anonymization within the civic dialogue network framework; that is, stories are conveyed in the same way as if listeners hear the original voice, and there is little to no difference in terms of the empathy or trust that listeners feel towards the storyteller.

Meanwhile, stories anonymized using TTS show significantly different scores than Original audio; there are statistically significant differences in 17/20 and 18/20 questions under the unaware and aware conditions, respectively. Notably, stories spoken using the TTS voice result in significantly lower scores in terms of perceived emotion (E1, E4), as well as trust, respect, and perception of authenticity (T1, T2, T3, T4, T5).

*5.1.2    Awareness of anonymization.* We now analyze the impact that awareness of anonymization has on listeners' empathy and trust for a storyteller. Fig. 2 shows scores for the Likert scale questions when performing head-to-head comparisons between the unaware and aware conditions for (a) Original, (b) VC, and (c) TTS voices. We again find largely similar patterns between the Original and VC voices. For the Original voice, there are no statistically significant differences between the two conditions for any of the 20 questions, while for the VC voice, there are statistically significant differences in only 2/20 questions. These results further support the notion that VC voices largely share the same characteristics as original, unaltered voices when perceived by listeners.

When looking at the scores for the TTS voice, however, we see an interesting phenomenon. When listeners are told that a TTS voice has been altered for the purposes of anonymization, they score the voice significantly higher in terms of perceived emotion (E1, E4), as well as trust, respect, and perceived authenticity (T1, T2, T3, T4, T5). Note that these are precisely the categories and questions in which TTS lagged behind the Original and VC voices in Section 5.1.1. This suggests that when people know that a voice has been altered, they become less sensitive to the vocal qualities of the storyteller, perhaps because they may guess that the voice was artifically synthesized (as was the case here). Consequently, the content of the story, rather than the delivery, may play a larger role in determining the listener's empathy and trust towards the storyteller.

*5.1.3    Free responses.* The analysis of the answers to the free response questions described in Section 4.3.2 yielded a set of 20 key topics and themes, which are listed in Table 3 along with the number of sentences from the responses belonging to each. These provide insights into the factors that participants in our survey were considering when they provided their answers to the Likert scale questions.

Participants generally cited positive sentiments when they provided high scores for empathy and trust, such as "genuine, real sounding experience" or "perceived passion and sincerity." One participant wrote, "I felt I could hear how passionate she was about her community, and heard her voice break when she spoke about her frustrations and disappointment that the state she lives in had found a new scapegoat." Another key factor seemed to be sympathy for the speaker that was elicited by their voice: "I could just feel the emotions in her voice, and although I may not have ever had any similar experiences in my life, I could feel considerable empathy for her." Notably, some participants mentioned that relatability of the story was a key factor for empathy, but trust could be determined more so by the voice and speaking style: "It wasn't a story that I could particularly relate to, but it seemed genuine and I had no reason to believe the speaker was inventing the story." Overall, these responses were most often connected to stories told using the Original or VC voices.

Meanwhile, lower empathy and trust scores were often linked to more neutral or negative sentiments, such as that the "voice sounded not genuine or unconvincing" or that the story "sounded scripted" or "like an AI." One participant responding to a TTS+unaware story stated, "judging from their voice, there wasn't much emotion there and I believe it was a mostly made up story," indicating that some listeners liken unnaturally un-emotive voices with insincerity or lower trust, regardless of the actual content of the story. However, other speakers were able to look past a lack of emotion and focus on the story's content: "I felt like the person's voice was pretty monotone and didn't convey emotion, but I do understand the situation the speaker is in and could identify with and empathize with that." Interestingly, the TTS+aware condition yielded some responses

Table 3. Topics and themes that emerged from answers to the free response question in Study 1 that asked, "What were some aspects of the audio clip that made you provide the answers that you did above?" The number of sentences from the free response answers belonging to each category are shown in parentheses.

| Topics and Themes | |
| --- | --- |
| 1. Genuine, real sounding experience (111) | 11. Relatable issue/content (83) |
| 2. Speaker sounds realistic and believable (16) | 12. Unrelatable issue/content (19) |
| 3. Strong emotion in the voice (106) | 13. Did not know enough about the speaker (4) |
| 4. Tone of voice (81) | 14. Voice sounded not genuine or unconvincing (20) |
| 5. Perceived passion and sincerity (46) | 15. Monotone voice, lack of emotion (94) |
| 6. Speaker spoke clearly and naturally (22) | 16. Skeptical that the experience happened (13) |
| 7. Powerful life experience (9) | 17. Story unclear or speaker difficult to understand (88) |
| 8. Sympathy for the speaker (18) | 18. Unclear audio (8) |
| 9. Sympathetic but unrelatable story (9) | 19. Sounds scripted / like reading from a script (30) |
| 10. Overall content of story (109) | 20. Doesn't sound like a real person / sounds like AI (42) |

like: "I could hear the emotion in her voice as she told her story, she sounded very genuine," and "The tone of voice and emotion when telling the story was authentic." These are sentiments that were rarely expressed in the TTS+unaware condition, and the responses possibly indicate that awareness of anonymization may have caused some listeners to subconsciously recalibrate their standards for expressivity and emotion.

## 5.2 Study 2

*5.2.1 Storytellers' perspectives on voice conversion vs. TTS for anonymization.* Fig. 3 shows a bar chart of average scores for the 5 Likert scale questions in the survey for Study 2, along with 95% confidence intervals. A clear pattern emerges when we analyze the perceptions of original storytellers towards stories that were anonymized using VC or TTS: VC voices consistently make speakers felt more heard, that the authenticity of their story would be retained, that the voice conveys the intended emotions, and that the anonymized voice retains important characteristics of their own voice. This is perhaps unsurprising given the additional prosodic information that VC preserves compared to TTS. However, the extent to which speakers show their preference is quite notable, as are the high scores that VC achieves. These results show that VC can successfully maintain the rich paralinguistic qualities of speech from the perspective of original speakers as well, fulfilling one of the key requirements for our civic dialogue network setting.

*5.2.2 Speakers' perceptions of privacy.* Participants rated both the VC and TTS voices highly for their perception of how successful they were at masking their identities; indeed, it was the only Likert scale question in Study 2 for which the two conditions did not have a statistically significant different score. Accordingly, when we looked at the community members with whom participants said they would be comfortable sharing their story with, we found that all participants marked exactly the same groups for versions of their story anonymized using both VC and TTS. Participants responded that they would be comfortable sharing their story in a VC or TTS voice with an average of $\mu = 4.95$ groups ($\sigma = 1.69$). This was a marked increase compared to the average number of groups that participants were comfortable sharing their story with in their original voice, which was $\mu = 2.95$ ($\sigma = 1.71$). Notably, 13 out of 21 participants stated that they would be comfortable sharing their story with all six groups mentioned in the survey (excluding "None of the above"), with a further two more (15 out of 21) stating that they would be comfortable sharing with every
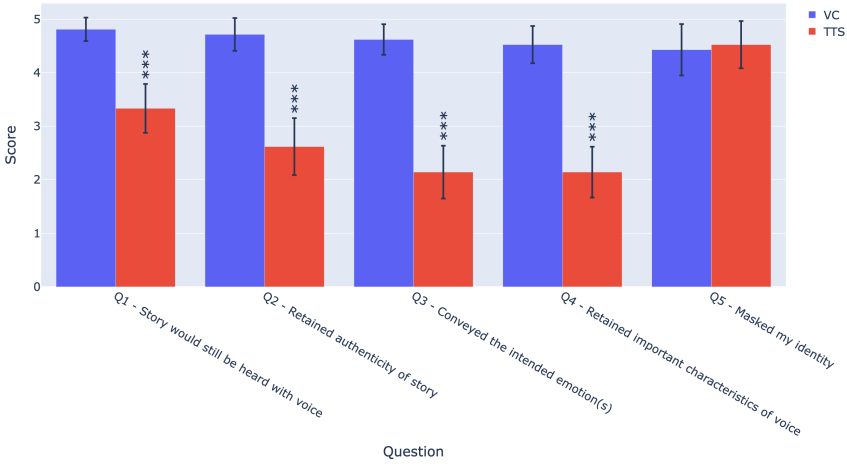
Fig. 3. Average scores for the 5 Likert scale questions in the survey for Study 2 comparing the perceptions of speakers towards their own voices that were anonymized using VC (blue) and TTS (red). Error bars represent 95% confidence intervals. *, **, or *** on top of bars for TTS denote statistically significant differences at $p < 0.05$, $p < 0.01$, or $p < 0.001$, respectively, for two-sided $t$-tests comparing against scores for VC.

group except for one. This was in contrast to the 3 out of 21 participants who were comfortable sharing their story with all six groups in their original voice.

## 6  Discussion

We have presented our findings from two studies: one measuring the impact of voice anonymization on listeners' feelings of empathy and trust towards individuals sharing personal stories, and one measuring the impact of voice anonymization on storytellers' own perceptions of being heard. Through these studies, we have shown that voice conversion is a technology that can be suitable for real settings that require anonymizing speech: it produces speech that has almost no difference with real voices in terms of listener perception, it preserves the emotion and prosody-related characteristics from the original speech, and most importantly, it can successfully mask a speaker's identity given an appropriate target speaker to convert the voice to. We have also shown that speech generated using TTS synthesis lags behind real and voice converted speech in these regards. However, our observation that awareness of speech modification significantly changes the way in which listeners perceive TTS voices opens up several interesting questions for how we might design speech interfaces that utilize this technology. For example, how might priming users of speech interfaces in different ways change how they react to the emotion and prosody, or lack thereof, expressed in TTS voices?

We additionally believe that our findings have meaningful implications for the design of technology-enhanced methods for civic dialogue. While voice is often used only as a metaphor for civic participation, there are cases in which spoken language can be a significant means of participation in civic systems, for example, as outlined in the civic dialogue network we considered in this work. Yet, mediums for richer communication often sacrifice opportunities for anonymity in exchange for the richness of the discourse [16]. In scenarios where spoken language is the primary communication medium, we show that anonymizing speech using voice conversion can be a legitimate and desirable way of allowing individuals to preserve richness and nuance in stories without sacrificing their ability to participate anonymously. Indeed, in numerous real-world applications of the civic

dialogue network, we have repeatedly observed requests for voice anonymization techniques that can fulfill this purpose, coming from communities ranging from company workplaces to academic institutions including high schools and universities. The widespread demand for these types of anonymization systems, whether in civic systems or beyond, can be further highlighted by the recent surge of work in this area, which has been spurred by initiatives such as the VoicePrivacy Challenge [76, 77].

Finally, we contribute to growing research within science, technology, and society studies as we explore not the effect of anonymity on behavior [10, 17, 19, 78], but its effect on trust and empathy from listeners and storytellers' perceptions of being heard. These characteristics are unique and important in public discourse, especially in today's civic era in which we see an increase in distrust and weak feelings of agency [36]. Given that anonymity can create conditions for discourse and deliberation free from coercion, our work takes a small step towards achieving the idealized vision of a liberal public sphere and a democracy filled with more equal, meaningful public discourse [32].

Computer-supported means of civic participation point to the desire for civic leaders and participants to have more nuanced, full opinions within their contributions [37, 44]. However, most explorations into anonymous participation are often held back by limited, survey-style evaluation standards. While we see some technologies such as Pol.is [1] experiment with anonymous discourse through chat, few look at anonymous discourse through voice, which further adds depth, richness, and nuance to one's civic contributions [16, 70]. We hope that our work may spark a discussion around more exploratory ways for engaging in civic technologies that invite vocal participation in pursuit of improving our public sphere.

## 7  Limitations

We finish by outlining some of the limitations of the present work. There are naturally scenarios in which neither voice conversion nor TTS are appropriate for performing anonymization. For example, certain voices in a social group may still be identifiable by their accent or idiolect if anonymization is done using voice conversion, since the technology as presented in this work entirely preserves a speaker's accent and speaking style. TTS offers an option for a different level of anonymization that could address some of these issues, and while we have seen that it removes the emotion and prosody-related characteristics from the original speech, it may still have a use case in situations in which voice conversion does not suffice. Leveraging other recent developments in speech synthesis technologies, such as accent conversion [38] or emotional TTS [79], may also help in this regard. On the other hand, changing speech characteristics such as accent or idiolect may affect listeners' perceptions of a story entirely. Any voice anonymization system should be chosen and implemented with care, considering the requirements and needs of the speakers, listeners, and stakeholders in the dialogue setting.

Social desirability bias, particularly in sensitive areas such as assessing one's warmth or empathy toward personal narratives, poses a challenge to the scientific integrity of studies like ours. To address this concern, we implemented several strategies to mitigate potential biases. First, we ensured randomization and uniformity across all survey conditions, thereby distributing any bias evenly across different stories. Additionally, participant responses were anonymized to promote candid feedback, and certain Likert-scale questions were reverse-coded to counteract response tendencies. However, it is important to acknowledge that our findings may still be influenced by an inherent bias toward higher empathy scores. Future research may address alternative methods for gauging empathy, such as eliciting responses based on how others might feel rather than solely focusing on participants' self-reported emotions.

In Study 2, conversations within a university setting were evaluated, involving faculty, staff, administrators, and students, encompassing diverse power dynamics and social complexities. However,

storytellers represented only a narrow segment of the broader, complex U.S. population, potentially limiting the study's generalizability. While Study 1 participants offered broader representation, their stories remained confined within a U.S. context as well, again potentially raising questions of generalizability. Future research should explore applicability across diverse cultural contexts and representative samples beyond the U.S. context.

## 8 Conclusion

In this work, we explored how to effectively incorporate a voice anonymization system into a spoken language-based civic dialogue network whose purpose is to increase participants' feelings of agency and being heard within civic processes. We tested two different speech transformation/synthesis methods for anonymization: voice conversion (VC) and text-to-speech (TTS). In contrast to previous work, we explored the effect of anonymity not on user behavior, but on empathy and trust elicited from listeners of stories, as well as on speakers' own perceptions of being heard. We found that voice conversion is a particularly suitable method for performing anonymization in our framework, as voice converted speech demonstrates almost no differences with real speech in terms of listeners' perceptions while successfully preserving the prosodic and other paralinguistic characteristics that were intended to be conveyed by the speaker. Our work contributes to literature on applications of speech-based technologies as well as on exploratory methods for engaging in discourse, and we hope that it may be a small step towards greater civic participation and improvement of our public sphere.

## References

[1] 2024. *Pol.is.* https://pol.is/home
[2] Dinislam Abdulgalimov, Reuben Kirkham, Stephen Lindsay, James Nicholson, Vasilis Vlachokyriakos, Emily Dao, Daniel Kos, Daniel Jitnah, Pam Briggs, and Patrick Olivier. 2023. Designing for the Embedding of Employee Voice. *Proceedings of the ACM on Human-Computer Interaction* 7, CSCW1, Article 45 (Apr 2023), 31 pages. https://doi.org/10.1145/3579478
[3] Philip E. Agre. 1999. Find Your Voice: Writing for a Webzine. https://pages.gseis.ucla.edu/faculty/agre/zine.html
[4] Danielle Allen. 2023. *Justice by Means of Democracy.* University of Chicago Press.
[5] Sherry R Arnstein. 1969. A ladder of citizen participation. *Journal of the American Institute of planners* 35, 4 (1969), 216–224.
[6] Hans Asenbaum. 2018. Anonymity and democracy: Absence as presence in the public sphere. *American Political Science Review* 112, 3 (2018), 459–472.
[7] Christina Baldwin and Ann Linnea. 2010. *The circle way: A leader in every chair.* Berrett-Koehler Publishers.
[8] Nancy K Baym. 2015. *Personal connections in the digital age.* John Wiley & Sons.
[9] Islam Bouzguenda, Chaham Alalouch, and Nadia Fava. 2019. Towards smart sustainable cities: A review of the role digital citizen participation could play in advancing social sustainability. *Sustainable Cities and Society* 50 (2019), 101627.
[10] Danah Boyd. 2012. The politics of "real names". *Commun. ACM* 55, 8 (2012), 29–31.
[11] Susan E Brennan and Maurice Williams. 1995. The feeling of another's knowing: Prosody and filled pauses as cues to listeners about the metacognitive states of speakers. *Journal of Memory and Language* 34, 3 (1995), 383–398.
[12] Julia Cambre, Jessica Colnago, Jim Maddock, Janice Tsai, and Jofish Kaye. 2020. Choice of voices: A large-scale evaluation of text-to-speech voice quality for long-form content. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems.* 1–13.
[13] Julia Cambre and Chinmay Kulkarni. 2019. One voice fits all? Social implications and research challenges of designing voices for smart devices. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 1–19.
[14] Edresson Casanova, Julian Weber, Christopher D Shulby, Arnaldo Candido Junior, Eren Gölge, and Moacir A Ponti. 2022. Yourtts: Towards zero-shot multi-speaker tts and zero-shot voice conversion for everyone. In *International Conference on Machine Learning.* PMLR, 2709–2720.
[15] Leigh Clark, Philip Doyle, Diego Garaialde, Emer Gilmartin, Stephan Schlögl, Jens Edlund, Matthew Aylett, João Cabral, Cosmin Munteanu, Justin Edwards, et al. 2019. The state of speech in HCI: Trends, themes and challenges. *Interacting with Computers* 31, 4 (2019), 349–371.
[16] Richard L Daft and Robert H Lengel. 1986. Organizational information requirements, media richness and structural design. *Management Science* 32, 5 (1986), 554–571.

[17]   Munmun De Choudhury and Sushovan De. 2014. Mental health discourse on reddit: Self-disclosure, social support, and anonymity. In *Proceedings of the International AAAI Conference on Web and Social Media*, Vol. 8. 71–80.

[18]   Tracey M Derwing, Marian J Rossiter, and Murray J Munro. 2002. Teaching native speakers to listen to foreign-accented speech. *Journal of multilingual and multicultural development* 23, 4 (2002), 245–259.

[19]   Judith S Donath. 2002. Identity and deception in the virtual community. In *Communities in Cyberspace*. Routledge, 37–68.

[20]   Pam Estell, Elizabeth Davidson, and Kaveh Abhari. 2021. Affording employee voice: how enterprise social networking sites (ESNS) create new pathways for employee expression. (2021).

[21]   Fuming Fang, Xin Wang, Junichi Yamagishi, Isao Echizen, Massimiliano Todisco, Nicholas Evans, and Jean-Francois Bonastre. 2019. Speaker Anonymization Using X-vector and Neural Waveform Models. In *10th ISCA Workshop on Speech Synthesis (SSW 10)*. ISCA.

[22]   Morris P Fiorina and Samuel J Abrams. 2008. Political polarization in the American public. *Annual Review of Political Science* 11 (2008), 563–588.

[23]   Patrick Fournier. 2011. *When citizens decide: Lessons from citizen assemblies on electoral reform.* Oxford University Press.

[24]   Nancy Fraser. 2014. Rethinking the public sphere: a contribution to the critique of actually existing democracy1. In *Between Borders*. Routledge, 74–98.

[25]   Archon Fung. 2006. Varieties of participation in complex governance. *Public Administration Review* 66 (2006), 66–75.

[26]   Akiko Fuse, Yuliya Navichkova, and Krysteena Alloggio. 2018. Perception of intelligibility and qualities of non-native accented speakers. *Journal of Communication Disorders* 71 (2018), 37–51.

[27]   Marshall Ganz. 2009. *Why David sometimes wins: Leadership, organization, and strategy in the California farm worker movement.* Oxford University Press.

[28]   Marshall Ganz, Nitin Nohria, and Rakesh Khurana. 2010. Leading change. In *Handbook of Leadership Theory and Practice: A Harvard Business School Centennial Colloquium.*

[29]   Weiyue Gao, Wei Xiang, Xuanhui Liu, Xueyou Wang, and Lingyun Sun. 2022. Impacts of Presenting Extra Information in Short Videos via Text and Voice on User Experience. In *2022 14th International Conference on Quality of Multimedia Experience (QoMEX)*. IEEE, 1–6.

[30]   Howard Giles. 1973. Communicative effectiveness as a function of accented speech. *Speech Monographs* 40, 4 (1973), 330–331. https://doi.org/10.1080/03637757309375813

[31]   Ioulia Grichkovtsova, Michel Morel, and Anne Lacheret. 2012. The role of voice quality and prosodic contour in affective speech perception. *Speech Communication* 54, 3 (2012), 414–429.

[32]   Jurgen Habermas. 1991. *The structural transformation of the public sphere: An inquiry into a category of bourgeois society.* MIT Press.

[33]   Yaowei Han, Sheng Li, Yang Cao, Qiang Ma, and Masatoshi Yoshikawa. 2020. Voice-indistinguishability: Protecting voiceprint in privacy-preserving speech data release. In *2020 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 1–6.

[34]   Daniel Hirst and Albert Di Cristo. 1998. A survey of intonation systems. *Intonation Systems: A Survey of Twenty Languages* 144 (1998), 152–166.

[35]   Chatham House. 2017. Chatham house rule.

[36]   Judith E Innes and David E Booher. 2004. Reframing public participation: strategies for the 21st century. *Planning Theory & Practice* 5, 4 (2004), 419–436.

[37]   Mahmood Jasim, Enamul Hoque, Ali Sarvghad, and Narges Mahyar. 2021. CommunityPulse: Facilitating Community Input Analysis by Surfacing Hidden Insights, Reflections, and Priorities. In *Proceedings of the 2021 ACM Designing Interactive Systems Conference* (Virtual Event, USA) *(DIS '21)*. Association for Computing Machinery, New York, NY, USA, 846–863. https://doi.org/10.1145/3461778.3462132

[38]   Mumin Jin, Prashant Serai, Jilong Wu, Andros Tjandra, Vimal Manohar, and Qing He. 2023. Voice-preserving zero-shot multiple accent conversion. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1–5.

[39]   Wonjune Kang, Mark Hasegawa-Johnson, and Deb Roy. 2023. End-to-End Zero-Shot Voice Conversion with Location-Variable Convolutions. In *Proc. INTERSPEECH 2023*. ISCA, 2303–2307.

[40]   Douglas Kellner. 2000. Habermas, the public sphere, and democracy: A critical intervention. *Perspectives on Habermas* 1, 1 (2000), 259–288.

[41]   Eugene Kharitonov, Damien Vincent, Zalán Borsos, Raphaël Marinier, Sertan Girgin, Olivier Pietquin, Matt Sharifi, Marco Tagliasacchi, and Neil Zeghidour. 2023. Speak, read and prompt: High-fidelity text-to-speech with minimal supervision. *arXiv preprint arXiv:2302.03540* (2023).

[42]   Jaehyeon Kim, Jungil Kong, and Juhee Son. 2021. Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech. In *International Conference on Machine Learning*. PMLR, 5530–5540.

[43] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020. Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis. *Advances in Neural Information Processing Systems* 33 (2020), 17022–17033.

[44] Narges Mahyar, Michael R. James, Michelle M. Ng, Reginald A. Wu, and Steven P. Dow. 2018. Community-Crit: Inviting the Public to Improve and Evaluate Urban Design Ideas through Micro-Activities. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (<conf-loc>, <city>Montreal QC</city>, <country>Canada</country>, </conf-loc>) *(CHI '18)*. Association for Computing Machinery, New York, NY, USA, 1–14. https://doi.org/10.1145/3173574.3173769

[45] Phil McAleer, Alexander Todorov, and Pascal Belin. 2014. How do you say 'Hello'? Personality impressions from brief novel voices. *PloS One* 9, 3 (2014), e90779.

[46] Siobhán McHugh. 2012. The Affective Power of Sound: Oral History on Radio. *Oral History Review* 39, 2 (2012).

[47] Leland McInnes, John Healy, and Steve Astels. 2017. hdbscan: Hierarchical density based clustering. *Journal of Open Source Software* 2, 11 (2017), 205.

[48] John Stuart Mill. 1966. *On liberty.* Springer.

[49] Raúl Montaño and Francesc Alías. 2016. The role of prosody and voice quality in indirect storytelling speech: Annotation methodology and expressive categories. *Speech Communication* 85 (2016), 8–18.

[50] Allison Morris and Gabrielle Maxwell. 2001. *Restorative justice for juveniles: Conferencing, mediation and circles.* Bloomsbury Publishing.

[51] Andreas Nautsch, Abelino Jiménez, Amos Treiber, Jascha Kolberg, Catherine Jasserand, Els Kindt, Héctor Delgado, Massimiliano Todisco, Mohamed Amine Hmani, Aymen Mtibaa, et al. 2019. Preserving privacy in speaker and speech characterisation. *Computer Speech & Language* 58 (2019), 441–480.

[52] Lilyana Ortega, Mikhail Lyubansky, Saundra Nettles, and Dorothy L Espelage. 2016. Outcomes of a restorative circles program in a high school setting. *Psychology of Violence* 6, 3 (2016), 459.

[53] Jose Patino, Natalia Tomashenko, Massimiliano Todisco, Andreas Nautsch, and Nicholas Evans. 2021. Speaker Anonymisation Using the McAdams Coefficient. In *Proc. INTERSPEECH 2021.* ISCA, 1099–1103.

[54] Pew Research Center. 2023. *Public Trust in Government: 1958-2023.* Technical Report. Washington, D.C. https://www.pewresearch.org/politics/2023/09/19/public-trust-in-government-1958-2023/

[55] Jeff Pittam. 1994. *Voice in Social Interaction: An Interdisciplinary Approach.* SAGE Publications.

[56] Miran Pobar and Ivo Ipšić. 2014. Online speaker de-identification using voice transformation. In *2014 37th International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO).* IEEE, 1264–1267.

[57] Jianwei Qian, Haohua Du, Jiahui Hou, Linlin Chen, Taeho Jung, Xiang-Yang Li, Yu Wang, and Yanbo Deng. 2017. Voicemask: Anonymize and sanitize voice input on mobile devices. *arXiv preprint arXiv:1711.11460* (2017).

[58] Kaizhi Qian, Yang Zhang, Shiyu Chang, Xuesong Yang, and Mark Hasegawa-Johnson. 2019. Autovc: Zero-shot voice style transfer with only autoencoder loss. In *International Conference on Machine Learning.* PMLR, 5210–5219.

[59] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning.* PMLR, 28492–28518.

[60] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP).* Association for Computational Linguistics.

[61] Howard Rheingold. 2008. *Using participatory media and public voice to encourage civic engagement.* MacArthur Foundation Digital Media and Learning Initiative.

[62] Emma Rodero. 2011. Intonation and emotion: influence of pitch levels and contour type on creating emotions. *Journal of Voice* 25, 1 (2011), e25–e34.

[63] Emma Rodero. 2015. The principle of distinctive and contrastive coherence of prosody in radio news: An analysis of perception and recognition. *Journal of Nonverbal Behavior* 39 (2015), 79–92.

[64] Emma Rodero. 2017. Effectiveness, attention, and recall of human and artificial voices in an advertising story. Prosody influence and functions of voices. *Computers in Human Behavior* 77 (2017), 336–346.

[65] Marcello Russo, Gazi Islam, and Burak Koyuncu. 2017. Non-native accents and stigma: How self-fulfilling prophesies can affect career outcomes. *Human Resource Management Review* 27, 3 (2017), 507–520.

[66] Angelien A Sanderman and René Collier. 1997. Prosodic phrasing and comprehension. *Language and Speech* 40, 4 (1997), 391–409.

[67] Klaus R Scherer, Rainer Banse, and Harald G Wallbott. 2001. Emotion inferences from vocal expression correlate across languages and cultures. *Journal of Cross-Cultural Psychology* 32, 1 (2001), 76–92.

[68] Ari Schlesinger, Eshwar Chandrasekharan, Christina A. Masden, Amy S. Bruckman, W. Keith Edwards, and Rebecca E. Grinter. 2017. Situated Anonymity: Impacts of Anonymity, Ephemerality, and Hyper-Locality on Social Media. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (Denver, Colorado, USA) *(CHI '17).* Association for Computing Machinery, New York, NY, USA, 6912–6924. https://doi.org/10.1145/3025453.3025682

[69] Juliana Schroeder and Nicholas Epley. 2015. The sound of intellect: Speech reveals a thoughtful mind, increasing a job candidate's appeal. *Psychological Science* 26, 6 (2015), 877–891.

[70] Juliana Schroeder, Michael Kardas, and Nicholas Epley. 2017. The humanizing voice: Speech reveals, and text conceals, a more thoughtful mind in the midst of disagreement. *Psychological Science* 28, 12 (2017), 1745–1762.

[71] Lijiang Shen. 2010. On a scale of state empathy during message processing. *Western Journal of Communication* 74, 5 (2010), 504–524.

[72] Berrak Sisman, Junichi Yamagishi, Simon King, and Haizhou Li. 2020. An overview of voice conversion and its challenges: From statistical modeling to deep learning. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 29 (2020), 132–157.

[73] Cass Sunstein. 1995. Democracy and the problem of free speech. *Publishing Research Quarterly* 11 (1995), 58–72.

[74] Deborah Tannen. 1982. *Spoken and Written Language: Exploring Orality and Literacy.* ABLEX Publishing Corporation.

[75] Alexis de Tocqueville. 2016. Democracy in america. In *Democracy: A Reader.* Columbia University Press, 67–76.

[76] Natalia Tomashenko, Xiaoxiao Miao, Pierre Champion, Sarina Meyer, Xin Wang, Emmanuel Vincent, Michele Panariello, Nicholas Evans, Junichi Yamagishi, and Massimiliano Todisco. 2024. The VoicePrivacy 2024 Challenge Evaluation Plan. *arXiv preprint arXiv:2404.02677* (2024).

[77] Natalia Tomashenko, Xin Wang, Xiaoxiao Miao, Hubert Nourtel, Pierre Champion, Massimiliano Todisco, Emmanuel Vincent, Nicholas Evans, Junichi Yamagishi, and Jean-François Bonastre. 2022. The VoicePrivacy 2022 Challenge Evaluation Plan. *arXiv preprint arXiv:2203.12468* (2022).

[78] Sherry Turkle. 2011. *Life on the Screen.* Simon and Schuster.

[79] Se-Yun Um, Sangshin Oh, Kyungguen Byun, Inseon Jang, ChungHyun Ahn, and Hong-Goo Kang. 2020. Emotional speech synthesis with rich and granularized control. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP).* IEEE, 7254–7258.

[80] Emily Van der Nagel and Jordan Frith. 2015. Anonymity, pseudonymity, and the agency of online identity: Examining the social practices of r/Gonewild. *First Monday* (2015).

[81] Junichi Yamagishi, Christophe Veaux, and Kirsten MacDonald. 2019. CSTR VCTK Corpus: English multi-speaker corpus for CSTR voice cloning toolkit (version 0.92). (2019).

[82] Guobin Yang. 2016. Narrative agency in hashtag activism: The case of #BlackLivesMatter. *Media and Communication* 4, 4 (2016), 13.