

---

# VoiceBox: Privacy through Real-Time Adversarial Attacks with Audio-to-Audio Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 As governments and corporations adopt deep learning systems to collect and analyze user-generated audio data, concerns about security and privacy naturally emerge in areas such as automatic speaker recognition. While audio adversarial examples offer one route to mislead or evade these invasive systems, they are typically crafted through time-intensive offline optimization, limiting their usefulness in streaming contexts. Inspired by architectures for audio-to-audio tasks such as denoising and speech enhancement, we propose a neural network model capable of adversarially modifying a user's audio stream in real-time. Our model learns to apply a time-varying finite impulse response (FIR) filter to outgoing audio, allowing for effective and inconspicuous perturbations on a small fixed delay suitable for streaming tasks. We demonstrate our model is highly effective at de-identifying user speech from speaker recognition and able to transfer to an unseen recognition system. We conduct a perceptual study and find that our method produces perturbations significantly less perceptible than baseline anonymization methods, when controlling for effectiveness. Finally, we provide an implementation of our model capable of running in real-time on a single CPU thread. Audio examples and code can be found at <https://master.d3hvhbnf7qxjtf.amplifyapp.com/>.

## 18 1 Introduction

19 Mass surveillance of voice communications is an ongoing and pervasive issue. While section 702 of the United States Foreign Intelligence Surveillance Act allows the government to perform targeted monitoring of foreign communications, bulk collection practices have resulted in the warrantless surveillance of large numbers of "incidental" foreign and domestic individuals [16]. Despite the fact that millions of these communications are obtained without warrants, they have been used in ordinary criminal investigations [58], undermining a core purpose of the Fourth Amendment of the US constitution [57]: to protect the people against searches without probable cause. Many corporations also possess the capability for large-scale collection of voice data [22], which may be leveraged to profile users for advertising or accessed by government entities through upstream and downstream surveillance [15]. As individuals are faced with these growing surveillance apparatus, it is worth remembering that routine surveillance of private voice communications is also corrosive to free speech and association and tends to disproportionately affect marginalized groups [27, 17].

31 In the absence of identifying metadata, automatic speaker recognition systems can facilitate mass surveillance by allowing an operator to search a database of recorded voice data for utterances from a chosen speaker [6] or to diarise (assign utterances to individuals) transcripts of recordings [43]. Prior to the advent of automatic speaker recognition these tasks required human analysts, forming a natural check on surveillance overreach. We seek to restore this check by degrading the efficacy of speaker recognition models while maintaining the original perceptual quality of the voice communication, a step that could grant users a measure of privacy from mass surveillance.

38 Modern speaker recognition systems rely on deep neural networks [5]. Deep networks have been  
39 shown to be vulnerable to adversarial examples—natural instances (e.g. a recording of "Bob"  
40 speaking) modified to cause a model to make an incorrect prediction (the recording’s speaker is  
41 labeled as "Maria") [52]. This presents opportunities for privacy-minded individuals to mislead or  
42 evade surveillance systems with adversarially-crafted inputs. Researchers have proposed adversarial  
43 attacks against a variety of audio systems, including speaker recognition [60, 50].

44 To fool a given neural network-based system, many audio attacks modify a recording by adding  
45 a perturbation signal directly at the waveform representation [7]. This perturbation is typically  
46 crafted using gradient information obtained from the system’s neural network—or from a similarly-  
47 constructed surrogate—and requires iterative optimization. Once the optimization is complete, the  
48 modified recording is played to the system in an effort to induce arbitrary incorrect predictions (an  
49 *untargeted* attack) or a specific incorrect prediction chosen by the attacker (a *targeted* attack).

50 In theory, such attacks might allow individuals to evade systems that surveil and analyze voice  
51 communications, thereby protecting privacy. However, many existing attack algorithms require a  
52 costly optimization for each audio recording to which they are applied. Continuous, real-time voice  
53 communication precludes the adoption of such approaches through online optimization (which is  
54 typically too slow) or through the use of a set of precomputed adversarial examples (which would  
55 necessarily limit the user’s interaction). This suggests a need for algorithms capable of modifying  
56 speech on-the-fly.

57 Recent years have seen the development of models for audio-to-audio tasks such as denoising, voice  
58 conversion, and musical timbre transfer [12, 47, 13]. Such models have sufficient expressive power to  
59 modify audio in coherent and task-specific ways, and are often designed to run in real-time. Inspired  
60 by these models, we propose **VoiceBox**, a deep network that learns to apply a time-varying finite  
61 impulse response (FIR) filter to outgoing audio, producing highly effective adversarial perturbations  
62 on a small fixed delay suitable for streaming. The main contributions of this work are:

- 63 • A highly effective method for de-identifying speech in real-time that is more perceptually  
64 inconspicuous than existing methods of equal effectiveness
- 65 • Objective and subjective experimental results that validate these claims
- 66 • A system that embodies the claimed advances, and that runs in real-time on a single CPU  
67 thread

68 Through this work, we hope to encourage further exploration of audio-to-audio models for protecting  
69 user privacy. Audio examples and code, including our streaming implementation, can be found at  
70 <https://master.d3hvhbnf7qxjtf.amplifyapp.com/>.

## 71 2 Related work

72 Previous works in **speech de-identification** have proposed methods for obfuscating speech to evade  
73 surveilling speaker-recognition systems. One such method is voice conversion, which modifies the  
74 speech of a source speaker to sound like that of another target speaker, both to humans and machines.  
75 Jin et al. [23] and Alegre et al. [4] analyze the vulnerability of speaker recognition systems against  
76 de-identification attacks performed with voice conversion models. However, the conversion models  
77 considered are incapable of real-time operation. While the recent availability of low-latency voice  
78 conversion models [47, 48, 30] may make such approaches more practical, our aim is to modify  
79 speech in a way that is inconspicuous and minimally invasive to the user experience, and that has the  
80 potential to generalize to tasks beyond speaker recognition. This also rules out recent anonymization  
81 approaches catalogued through the VoicePrivacy Challenge, as all evaluated submissions to date  
82 noticeably alter speaker characteristics [54]. Similarly, we do not consider obvious transformations  
83 like pitch-shifting that significantly alter or degrade recorded audio.

84 **Real-time audio attacks** have recently been explored. Chiquier et al. [9] propose an attack on  
85 automatic speech recognition systems that is capable of inducing significant transcription errors in an  
86 *over-the-air* environment (where attack audio is played to the system in a physical space), but which  
87 introduces conspicuous noise due to the use of additive perturbations at the audio waveform. Their  
88 approach also necessitates both a large model (requiring approximately 16 GPU-days to train) and  
89 lengthy receptive field, resulting in an initial “idle” period of 2.5s during which the attack does not

90 affect speech. By contrast, we focus on an *over-the-line* setting, where attacks are passed to the victim  
91 model over a digital channel, and propose a lightweight network capable of producing inconspicuous  
92 perturbations through the use of filtering, with an initial delay of miliseconds rather than seconds.

93 Universal adversarial attacks optimize a short perturbation that can be played alongside arbitrary  
94 speech, optionally in real-time, to evade a victim model [31, 60, 26]. However, these attacks tend to  
95 introduce conspicuous noise due to the use of additive perturbations at the audio waveform. The same  
96 holds for the attack of Xie et al. [59], which uses a Wave-U-Net [51] model to efficiently generate  
97 universal perturbations. Rather than generate additive perturbations, we perform multiplicative attacks  
98 in the frequency domain through the use of filtering, and in doing so avoid the bias towards noise-like  
99 artifacts typical of additive attacks.

100 Other works have proposed crafting **adversarial attacks in the frequency domain**. The “Ke-  
101 nansville” attack of Abdullah et al. [1] performs spectral gating, removing low-energy frequency  
102 components in an effort to hinder classification. The authors demonstrate the effectiveness of their  
103 approach against seen and unseen speaker- and speech-recognition systems, and the proposed method  
104 is fast and gradient-free. However, the spectral gating must be constrained to a small number of  
105 key audio frames to avoid conspicuous artifacts, which requires offline optimization using word- or  
106 phoneme-level alignment information.

107 Ahmed et al. [3] optimize bandpass filters in order to perform targeted impersonation attacks on  
108 speaker-recognition models; these filters can then be physically realized as resonant tubes through  
109 which an attacker can speak to the victim system, allowing real-time operation. However, once  
110 optimized and realized, the bandpass filters are fixed in time and cannot adapt to the attacker’s speech.  
111 The attack is significantly less effective than traditional approaches, even in controlled acoustic  
112 environments, and requires an active intervention on the part of the user to both realize and apply.

113 Finally, O’Reilly et al. [40] adversarially optimize the parameters of a time-varying finite impulse  
114 response (FIR) filter in order to avoid noisy artifacts. However, the attack operates offline and must  
115 be optimized separately for every instance to which it is applied.

116 In summary, we are aware of no current approach that is simultaneously capable of performing de-  
117 identification while remaining inconspicuous to human listeners and operating in real-time on-device.

### 118 3 VoiceBox

119 We propose VoiceBox, a system that applies adversarial time-varying filtering to user audio in  
120 real-time. VoiceBox modifies speech by specifying the frequency magnitude response of a standard  
121 finite impulse response (FIR) filter. By carefully varying the filter’s response many times per second,  
122 the speech is de-identified to an automated system while speaker identity is preserved for human  
123 listeners. This is accomplished without introducing noisy artifacts at the waveform, as filtering can  
124 only amplify or attenuate frequency energy present in the original signal through multiplication in the  
125 Fourier domain (see Section 3.3 for details of the filtering implementation used in our experiments).

126 VoiceBox consists of three main modules, shown in Figure 1: (1) an **encoder** module extracts  
127 acoustic features from input audio frames, (2) a recurrent **bottleneck** module incorporates context  
128 from past frames to predict a set of adversarial filter controls for each frame, and (3) a **decoder**  
129 module regularizes the predicted filter controls and applies them to the corresponding frame of input  
130 audio. We provide an overview of each component below, and further details in Appendix A.

#### 131 3.1 Encoder

132 The encoder module extracts acoustic features from input audio to guide adversarial filtering. Our  
133 network accepts 16kHz audio (double the sample rate of telephone speech) segmented into frames of  
134 256 samples with 50% overlap. Though our model is capable of operating in fully causal fashion, we  
135 find that we achieve stronger attacks using a lookahead of five frames (see Section 3.2), resulting in a  
136 minimum theoretical latency of 56ms.

137 Taking cues from voice-conversion systems, we aim to disentangle speaker characteristics from the  
138 linguistic content of input utterances so that our model can subtly manipulate the former. Audio  
139 frames entering the encoder are passed to four sub-modules in parallel to obtain the following features.

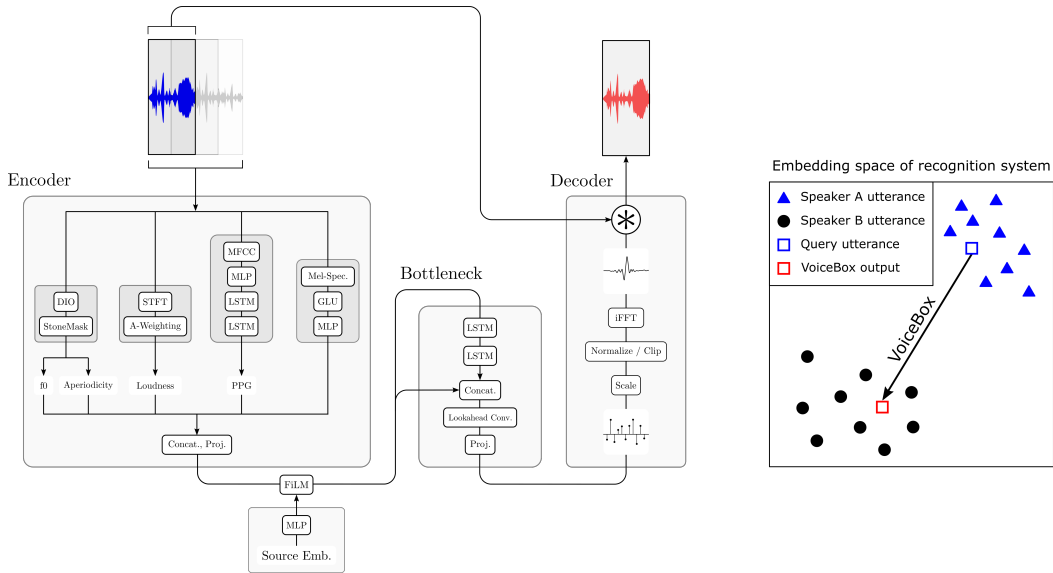


Figure 1: **Left:** The proposed VoiceBox architecture. Acoustic features extracted by the encoder are fed to the recurrent bottleneck to predict filtering controls, which are regularized and applied to the input by the decoder to obtain adversarial audio. **Right:** VoiceBox adversarially perturbs the user’s audio stream such that any extracted queries are scored by the system as dissimilar to the user’s enrolled utterances, hampering identification or retrieval.

140 **Pitch features:** Given the known sensitivity of speaker-recognition models to pitch [3] and the  
 141 importance of spectral structure in delineating linguistic content, we extract fundamental frequency  
 142 and aperiodicity estimates for each frame using the DIO algorithm [35]. Pitch estimates are refined  
 143 by the StoneMask algorithm [36]. For both stages, we use the `pyworld` [21] implementation of the  
 144 WORLD vocoder [36] (MIT license).

145 **Loudness:** We take an A-weighted average of each frame’s log-magnitude spectrum to obtain a  
 146 loudness estimate [34].

147 **Phonetic posteriorgrams:** Following the method of Ronssin & Cernak [47], we use a trained  
 148 phoneme classifier with frozen weights to encode linguistic content. The classifier consists of a  
 149 multi-layer perceptron followed by two LSTM [20] layers and a linear classification layer, and takes  
 150 as input 13 mel-frequency cepstral coefficients (MFCC) with first- and second-order deltas for each  
 151 frame. We train our phoneme classifier on the "train-clean-100" subset of the LibriSpeech dataset  
 152 [41] using frame-aligned phoneme labels [32, 33] (CC-BY 4.0 license). Rather than directly use the  
 153 classifier’s predicted distributions over phoneme labels – known as phonetic posteriorgrams (PPGs)  
 154 [18] – we discard the classification layer after training and pass along the output of the final LSTM  
 155 layer, which Ronssin & Cernak also refer to as PPGs in their architecture.

156 **Spectrogram features:** Finally, we use a simple network consisting of a gated linear unit and  
 157 multi-layer perceptron operating independently on each mel-spectrogram frame—a widely used  
 158 speech representation [25]—to capture any residual information.

159 For each frame, the above features are concatenated and projected linearly to obtain a low-dimensional  
 160 encoding. We then introduce additional speaker information by passing a fixed, pre-computed  
 161 embedding of the source speaker through a FiLM layer [44] to modulate the encoder output. This  
 162 embedding helps to guide the de-identification task, allowing us to train a single model capable of  
 163 de-identifying arbitrary users. This imposes a requirement that users must record a small amount of  
 164 speech to obtain an embedding before first using VoiceBox. We find that less than a minute of speech  
 165 is sufficient, and use a pre-trained ResNetSE34V2 model [19] to compute embeddings (see sections  
 166 4.1, 4.2).

### 167 3.2 Bottleneck

168 Encodings entering the bottleneck are passed through two LSTM layers, and the outputs are con-  
169 catenated with a skip connection from the encoder. To enable streaming, we use unidirectional  
170 LSTM layers and pass concatenated outputs through a small lookahead convolutional network [55] to  
171 incorporate information from future frames at the expense of a small fixed delay. We find a lookahead  
172 of 5 frames (roughly 48 ms at our sample rate) is sufficient to craft strong de-identification attacks.  
173 Finally, a linear projection layer maps the concatenated representations to a vector of filter controls,  
174 with each entry representing the unnormalized frequency magnitude response of a filter band for the  
175 current frame.

### 176 3.3 Decoder

177 Our decoder applies time-varying filtering to the input audio based on frame-wise controls obtained  
178 from the bottleneck. For our task, we find that a filtering-based decoder affords a number of  
179 advantages over the synthesis-based (e.g. transposed-convolutional) decoder architectures present in  
180 many audio-to-audio models [12]:

- 181 • The filtering module is not prone to the periodic upsampling artifacts that transposed-  
182 convolutional architectures often introduce [45].
- 183 • Our decoder is a simple deterministic module with no trainable parameters, keeping Voice-  
184 Box lightweight.
- 185 • The capacity of the decoder and conspicuousness of perturbations can easily be constrained  
186 in terms of the number of filter bands or their allowed range of motion, providing inter-  
187 pretable control to the user.

188 To regularize and apply filter controls to each frame of audio, we use a method similar to that of  
189 O’Reilly et al. and Engel et al. [40, 13]. We apply sigmoid scaling to bound filter controls to the  
190 range  $[0, 2]$ . We clip deviations from unity beyond a fixed bound  $\epsilon$ . Each set of scaled filter controls  
191 is transformed into a time-domain impulse response via the inverse Fourier transform. We shift the  
192 impulse response to zero-phase (symmetric) form, apply a Hann window, and finally convolve with  
193 the corresponding input audio frame by taking the Fourier transform and performing element-wise  
194 multiplication. After compensating for shift, we overlap-add the resulting frames with a Hann window  
195 to obtain the final filtered audio. We discuss the details of our buffered streaming implementation in  
196 A.3. Our VoiceBox model is implemented in PyTorch [2] and contains 6.3m trainable parameters,  
197 and 7.5m in total counting the frozen phoneme encoder.

### 198 3.4 Training objective

199 To demonstrate the ability of VoiceBox to perform inconspicuous privacy-preserving audio trans-  
200 formations, we train our model to attack speaker recognition systems. Given a large database of  
201 speech recordings and access to a user’s audio stream, a surveilling entity may seek to (a) identify  
202 the user by matching their speech against recordings of known provenance in the database, or (b)  
203 retrieve other utterances of the user from the database. Systems designed for these tasks – speaker  
204 recognition and retrieval, respectively – often rely on neural network models to map speech utterances  
205 to a low-dimensional embedding space in which distance corresponds to speaker similarity [14]. In  
206 this work, we consider models  $f$  for which the speaker distance  $D_f$  between utterances  $u$  and  $v$  can  
207 be measured via a cosine distance between embeddings  $f(u)$  and  $f(v)$ :

$$D_f(u, v) = 1 - \frac{f(u) \cdot f(v)}{\|f(u)\|_2 \|f(v)\|_2} \quad (1)$$

208 At inference time, *query* audio from the user is embedded and scored for similarity against *enrolled*  
209 utterances – pre-computed embeddings stored in the database. Scores may be evaluated for all  
210 enrolled embeddings, or for only a representative of each unique speaker (e.g. speaker centroids).  
211 The highest-ranking result or results (e.g. the closest speaker identity) are then returned. We do not  
212 distinguish between recognition and retrieval tasks, and refer to both under the umbrella of "speaker  
213 recognition." This is because from a privacy standpoint, the objective in each task is identical: alter  
214 the query audio to prevent valid matches, thereby de-identifying the user.

215 A variety of methods have been proposed to efficiently compute and search over low-dimensional  
 216 speaker representations. Generally, a hashing algorithm is applied to speaker embedding vectors  
 217 to reduce storage and search costs [49, 29, 14]. Because the resulting hash representation merely  
 218 serves as an efficient point of access for embedding-space distances, we omit hashing algorithms  
 219 from consideration and define our attack objectives on the embedding space directly.

220 Given a speaker embedding model, we aim to modify query audio drawn from a user’s stream  
 221 such that its embedding-space distance from any enrolled utterances of the user is large, and its  
 222 evaluated similarity is small. We quantify this de-identification in terms of distance thresholds in  
 223 the embedding space, set according to percentiles of the estimated distribution of all inter-speaker  
 224 embedding distances. Let  $P_r$  represent the distance corresponding to the  $r^{\text{th}}$  percentile of this  
 225 distribution; then for query audio  $u$  and enrolled utterance  $v$ ,  $D_f(u, v) > P_r$  implies that roughly  $r$   
 226 percent of database entries should be scored as more similar to  $u$  than  $v$ . Thus, we can set distance  
 227 thresholds that correspond directly to the strength of de-identification applied to user audio, and  
 228 construct a loss function that penalizes our model when the embedding-space cosine distance between  
 229 query and enrolled utterances falls below the threshold. To do so, we use a variant of the adversarial  
 230 loss proposed by Zhang et al. [60]. Let  $f$  represent the victim model,  $g$  our VoiceBox network,  $u$  an  
 231 utterance from our user’s audio stream, and  $P_r$  our de-identification threshold; then

$$\mathcal{L}_{adv}(f, g, u) = (P_r - D_f(g(u), u) + \kappa)^+ \quad (2)$$

232 where  $(\cdot)^+ = \max(\cdot, 0)$  and  $\kappa$  is a confidence parameter encouraging the attack to fully cross the  
 233 threshold. To ensure that our VoiceBox model learns to perturb user audio inconspicuously, we  
 234 incorporate an additional loss function to penalize perceptible filtering artifacts. We compute the  
 235 combined waveform  $L_1$  and multi-resolution spectrogram losses proposed by Defosséz et al. [12] on  
 236 the clean and adversarial audio:

$$\mathcal{L}_{aux}(g, u) = \mathcal{L}_{stft}(u, g(u)) + \|u - g(u)\|_1 \quad (3)$$

237 where  $\mathcal{L}_{stft}$  is given by a sum of magnitude and spectral convergence losses computed over sev-  
 238 eral spectrogram resolutions. We provide further details in Appendix A.2. Combining the above  
 239 adversarial and auxiliary losses, we obtain our final attack objective:

$$\mathcal{L} = \mathcal{L}_{adv} + \mathcal{L}_{aux} \quad (4)$$

## 240 4 Experimental design

241 We describe experiments used to validate the claimed advances of our work, namely that VoiceBox  
 242 can de-identify speech in real-time while remaining significantly less conspicuous than existing  
 243 methods of similar effectiveness. We first introduce the models, datasets, and attacks considered in  
 244 our experiments. Following this, we detail our experiment configurations and present the results of  
 245 both objective and subjective evaluations.

### 246 4.1 Speaker recognition models

247 **ResNetSE34v2:** We train attacks against the ResNetSE34v2 model [19] provided in the  
 248 VoxCeleb Trainer repository [11] (MIT License). The model takes mel-spectrogram inputs and  
 249 uses 2D convolutions with residual connections, squeeze-and-excitation, and attentive statistics pool-  
 250 ing to generate frame-level features and aggregate them into 512-dimensional speaker embeddings.

251 **Y-Vector:** To examine the transferability of our approach against unseen systems, we evaluate  
 252 trained attacks against the Y-Vector model proposed by Zhu et al. [61]. The model uses a multiscale  
 253 1D-convolutional waveform encoder to extract acoustic features, followed by squeeze-and-excitation  
 254 blocks, feature aggregation, and a time-delayed neural network to map variable-length utterances to  
 255 128-dimensional speaker embeddings.

256 Both the ResNetSE34v2 and Y-Vector models were trained on the development set of the VoxCeleb2  
 257 dataset [10]. Note that while our attack is causal—modulo a short, five-frame lookahead—we do not  
 258 impose the same restriction on a hypothetical surveillance system; instead, we evaluate our attack  
 259 against strong, non-causal systems capable of aggregating speaker characteristics from across full  
 260 utterances before rendering predictions.

## 261 4.2 Datasets

262 **LibriSpeech:** (CC-BY 4.0) We use both the `train-clean-100` and `test-clean` subsets of the  
263 LibriSpeech dataset [41] for training VoiceBox. The former comprises 28,539 utterances from 251  
264 speakers while the latter comprises 2,620 utterances from 40 speakers. We trim or pad all utterances  
265 to 4 seconds.

266 **VoxCeleb1:** (CC-BY 4.0) To simulate large-scale surveiling speaker recognition, we evaluate attacks  
267 on the VoxCeleb1 dataset [38], comprising 153,516 utterances from 1,251 speakers. This also ensures  
268 that no evaluation speakers are seen during training. As with the LibriSpeech dataset, we trim or pad  
269 all utterances to 4 seconds.

270 For all experiments, we carefully divide the data to imitate a realistic attack setting. During training,  
271 we select fifteen utterances (one minute total) from each source speaker in the training set and compute  
272 embeddings using the ResNetSE34v2 (MIT License) model. The centroid of these embeddings is  
273 then used as an enrolled target for the computation of the adversarial loss with all utterances of that  
274 speaker. We find that this produces stronger attacks than using individual utterance embeddings as  
275 targets, possibly by ensuring more consistent gradient information across the optimization. For our  
276 VoiceBox attack (see Section 4.3), we select a further ten utterances (40s total) from each source  
277 speaker in the training set and again compute embeddings using the ResNetSE34v2 model. The  
278 centroid of these embeddings is then fed as a fixed conditioning vector to the VoiceBox model  
279 alongside all utterances from that speaker (see Section 3.1). Similar to training, during evaluation we  
280 select fifteen utterances per speaker as a query set. We again select a further ten utterances to serve as  
281 conditioning for the VoiceBox attack. Finally, twenty utterances of each speaker are enrolled in the  
282 speaker recognition system, serving as the database against which query utterances are matched.

## 283 4.3 Attack algorithms

284 We perform untargeted attacks using the following algorithms. **VoiceBox:** We implement the  
285 proposed VoiceBox attack as described in Section 3 and Appendix A and train for 10 epochs on 3  
286 NVidia RTX 2080 Ti GPUs; this takes approximately 40 minutes. **Universal:** We optimize a short  
287 (2s) universal additive perturbation for 10 epochs using the established penalty method [31, 60, 26]  
288 and the same adversarial objective as VoiceBox. On 3 NVidia RTX 2080 Ti GPUs, this takes  
289 approximately 16 minutes. To limit the perceptibility of the attack, we scale the perturbation to have  
290  $L_\infty$  norm 0.08 times that of the unperturbed speech. During training and evaluation, the perturbation  
291 is aligned arbitrarily with query utterances and looped to match durations, serving as a constant  
292 adversarial "background" signal. **White noise:** We add Gaussian noise to utterances at the waveform  
293 representation at a signal-to-noise ratio of  $-10$ dB. **Spectral gating:** We modify the "Kenansville"  
294 attack of Abdullah et al. [1] to allow for streaming use by performing spectral gating at all frames,  
295 using a threshold of 4dB relative to the maximum-energy spectral bin of each frame.

## 296 4.4 Objective evaluation

297 We evaluate attacks in a large-scale closed-set speaker recognition task. First, we use the `test-clean`  
298 LibriSpeech subset to obtain a rough estimate of the distribution of distances between an individual  
299 utterance and the centroids of each distinct speaker in the embedding space of the ResNetSE34V2  
300 model. To encourage strong de-identification attacks, we take the distance corresponding to the 25<sup>th</sup>  
301 percentile of this distribution as the target threshold  $P_{25}$  for our training loss (see Section 3.4). We  
302 train each attack as discussed above.

303 We evaluate the closed-set recognition of all attacks over the VoxCeleb1 dataset. We compute  
304 the distance between each query embedding and the centroid of all embeddings of each speaker;  
305 recognition is then performed by returning the speaker identity of the nearest speaker centroid. We  
306 find this is slightly more accurate and robust to attack than using the identity of the nearest embedded  
307 utterance. We perform exact nearest-neighbors search over the embedding space using FAISS  
308 [24] (MIT license), and report the top-1 (T-1) and top-10 (T-10) accuracy of the relevant speaker  
309 recognition model given both clean and adversarial queries. Additionally, we compute the following  
310 objective speech quality metrics over the clean and adversarial query audio as a proxy measure of  
311 the imperceptibility of attacks: Perceptual Evaluation of Speech Quality (PESQ) [46], operating  
312 in the wide-band configuration and using the `python-pesq` implementation [56] (MIT license);

Table 1: Results of our objective evaluation. We perform attacks on the VoxCeleb1 dataset against seen (ResNetSe34v2) and unseen (Y-Vector) speaker recognition models, and compute both top-1 and top-10 recognition accuracies. Additionally, we compute a set of objective speech quality metrics for each attack.

Approach	Speech Quality Metrics		ResNetSe34V2		Y-Vector	
	PESQ $\uparrow$	STOI $\uparrow$	T-1 $\downarrow$	T-10 $\downarrow$	T-1 $\downarrow$	T-10 $\downarrow$
White noise	1.03	0.23	0.13	0.40	0.00	0.01
Spectral gating	1.12	0.56	0.02	0.11	0.02	0.12
Universal	2.36	0.82	0.14	0.22	0.44	0.64
VoiceBox	3.77	0.90	0.02	0.10	0.32	0.62
No attack	4.64	0.99	0.97	0.99	0.93	0.98

313 and Short-Time Objective Intelligibility (STOI) [53], using the `pystoi` implementation [42] (MIT  
314 license). The results of our evaluation are presented in table 1.

315 We evaluate the real-time performance of the streaming implementation of the VoiceBox attack  
316 described in Appendix A.3 by measuring its average real-time factor (RTF). We measure performance  
317 on a single thread on two different CPUs, an Intel i7-5600U @ 3.2 GHz and an Apple M1 Chip. In the  
318 streaming configuration, VoiceBox processes a chunk of 4 overlapping frames at a time, equivalent to  
319 640 samples / 40 ms. We compute the average RTF for processing a chunk over all chunks in a 4  
320 second audio clip, and find that VoiceBox has an average RTF of .255 on the Intel i7-5600U and .200  
321 on the Apple M1.

#### 322 4.5 Subjective evaluation

323 For our subjective evaluation, we sampled 100 clean speech recordings from the VoxCeleb1 dataset.  
324 Each one of the four attacks described in Section 4.3 was applied to all 100 recordings, resulting in  
325 100 comparison sets with five recordings per set: the original clean speech and the speech modified  
326 by each of the four attacks. These sets were then evaluated in a MUSHRA-style listening study [8]  
327 deployed on Amazon Mechanical Turk using the open-source, MIT-licensed Reproducible Subjective  
328 Evaluation (ReSEval) [37] system. IRB approval was obtained prior to conducting this study, and  
329 there are no known risks to the participants in this study.

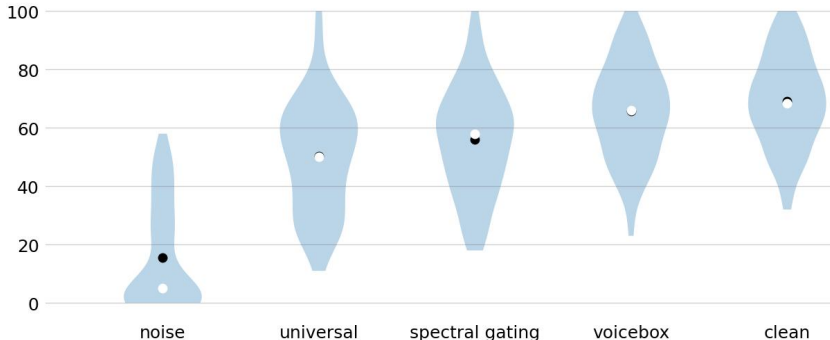


Figure 2: Distributions of quality ratings from our crowdsourced subjective listening MUSHRA-style test on audio quality. Higher numbers are better. Black dots are means and white dots are medians. Wilcoxon signed-ranked tests between all pairs of conditions show statistical significance at  $p < 0.05$ .

330 We recruited 20 participants. Participants were screened with a listening test prior to beginning the  
331 study. Each participant that passed the screening rated 20 comparison sets. For each comparison set,  
332 the participant was asked to listen to and rate the relative audio quality of each of the five audio files  
333 on a scale from 0 to 100. We omitted responses by four participants who failed our prescreening  
334 listening test and nine who rated the white noise attack (our low anchor) as superior in quality to  
335 ground-truth clean audio (our high anchor), giving us a total of 140 five-way comparisons. For more  
336 details of our crowdsourced subjective evaluation, see Appendix C.

337 The results of our crowdsourced subjective evaluation can be found in Figure 2. We find VoiceBox  
338 is preferred to all other methods, and exhibits similar perceptual quality to ground-truth speech  
339 recordings. The difference between VoiceBox and ground-truth audio is significant ( $p < 0.05$ ) using  
340 a Wilcoxon signed-rank test ( $p = 0.024$ ) but not significant using Welch’s T-test ( $p = 0.071$ ), which  
341 assumes that human perceptual scores are normally distributed but with potentially different variances.  
342 All other pairs of conditions are significantly different using either test.

#### 343 4.6 Additional experiments

344 We conduct supplementary experiments demonstrating the robustness of VoiceBox against a deep  
345 network-based speech enhancement model and examining its performance under the assumption that  
346 adversarial queries are enrolled by the surveiling speaker recognition system. These experiments are  
347 detailed in appendices B.1 and B.2, respectively.

### 348 5 Discussion

349 Our experimental results indicate VoiceBox can de-identify speech from arbitrary unseen users in  
350 real-time on a standard M1 CPU. The de-identified speech is of significantly higher audio quality  
351 than competing methods, as reported by a subjective listener study and as measured through standard  
352 metrics of speech intelligibility. Our method also achieves nontrivial de-identification results against  
353 a system it was not trained on, outperforming a far more perceptible universal attack. This is  
354 notable given that we take no explicit steps to improve the transferability of our method. While  
355 pure signal-processing approaches such as spectral gating and white noise transfer between systems  
356 more successfully, they severely degrade audio quality, resulting in word-error-rates 10 to 20 times  
357 higher than VoiceBox and much lower listener quality ratings. This limits their practicality in real-  
358 world voice communications. By contrast, our method produces adversarial audio that is virtually  
359 indistinguishable from the clean source, as indicated by our subjective evaluation.

360 We believe the benefits of real-time, imperceptible de-identification attacks for user privacy are  
361 self-evident. Our same approach, however, could be argued to hamper legitimate targeted surveillance  
362 (e.g. tracking of the communications of criminal organizations). This ability to evade tracking  
363 is, however, limited by the fact that VoiceBox’s output is perceptually quite close to the original  
364 speech. Therefore, a human listener would still be able to readily identify the speaker. As such, the  
365 effect of VoiceBox would simply be to return us to the status quo that existed prior to automated  
366 large-scale speaker recognition of requiring manual human evaluation. VoiceBox could also, in  
367 concept, be used to fraudulently access information or services protected by speaker verification.  
368 However, since VoiceBox performs perceptually inconspicuous filtering, the speech produced would  
369 still sound like the original speaker and could not pass human inspection. Both voice-conversion  
370 and speech synthesis-based attacks produce speech that is perceptually similar to the target and  
371 are, therefore, much more suited to this purpose. Our approach can be thought of as leveraging the  
372 asymmetry between system and attacker attention inherent to mass surveillance, using inconspicuous  
373 perturbations to evade large-scale systems that must render predictions in bulk. This same asymmetry  
374 is not necessarily present when trying to bypass authentication mechanisms and access tightly-guarded  
375 services.

376 We view our method as an initial step towards protecting users from indiscriminate mass surveillance.  
377 A number of obvious directions for future work stand out, such as improving the transferability of  
378 adversarial examples [39], and evaluating attacks against real-world speaker recognition systems and  
379 over real-world communication channels. We hope this work encourages further exploration of the  
380 applications of audio-to-audio models for protecting user privacy.

### 381 References

- 382 [1] Hadi Abdullah, Muhammad Rahman, Washington Garcia, Kevin Warren, Anurag Yadav, Tom  
383 Shrimpton, and Patrick Traynor. Hear "no evil", see "kenansville"\*: Efficient and transferable  
384 black-box attacks on speech recognition and voice identification systems. In *IEEE Symposium*  
385 *on Security and Privacy*, 2021.
- 386 [2] Francisco Massa Adam Lerer James Bradbury Gregory Chanan Trevor Killeen Zeming Lin  
387 Natalia Gimelshein Luca Antiga Alban Desmaison Andreas Kopf Edward Yang Zachary

- 388 DeVito Martin Raison Alykhan Tejani Sasank Chilamkurthy Benoit Steiner Lu Fang Jun-  
389 jie Bai Adam Paszke, Sam Gross and Soumith Chintala. Pytorch: An imperative style, high-  
390 performance deep learning library. In *Neural Information Processing Systems (NeurIPS)*,  
391 2019.
- 392 [3] Shima Ahmed, Yash Wani, Ali Shahin Shamsabadi, Mohammad Yaghini, Iliia Shumailov,  
393 Nicolas Papernot, and Kassem Fawaz. Pipe overflow: Smashing voice authentication for fun  
394 and profit. *arXiv preprint arXiv:2107.14642*, 2022.
- 395 [4] Federico Alegre, Giovanni Soldi, Nicholas Evans, Benoit Fauve, and Jasmin Liu. Evasion and  
396 obfuscation in speaker recognition surveillance and forensics. In *International Workshop on*  
397 *Biometrics and Forensics*, 2014.
- 398 [5] Zhongxin Bai and Xiao-Lei Zhang. Speaker recognition based on deep learning: An overview.  
399 *Neural Networks*, 140:65–99, 2021.
- 400 [6] Peter J. Barger and Sridha Sridharan. On the performance and use of speaker recognition  
401 systems for surveillance. In *International Conference on Advanced Video and Signal Based*  
402 *Surveillance (AVSS)*, 2006.
- 403 [7] Nicholas Carlini and David Wagner. Audio adversarial examples: Targeted attacks on speech-  
404 to-text. In *IEEE Security and Privacy Workshops*, 2018.
- 405 [8] Mark Cartwright, Bryan Pardo, Gautham J Mysore, and Matt Hoffman. Fast and easy crowd-  
406 sourced perceptual audio evaluation. In *International Conference on Acoustics, Speech and*  
407 *Signal Processing (ICASSP)*, 2016.
- 408 [9] Mia Chiquier, Chengzhi Mao, and Carl Vondrick. Real-time neural voice camouflage. In  
409 *International Conference on Learning Representations (ICLR)*, 2022.
- 410 [10] J. S. Chung, A. Nagrani, and A. Zisserman. Voxceleb2: Deep speaker recognition. In *Inter-*  
411 *speech*, 2018.
- 412 [11] Joon Son Chung, Jaesung Huh, Seongkyu Mun, Minjae Lee, Hee Soo Heo, Soyeon Choe,  
413 Chiheon Ham, Sunghwan Jung, Bong-Jin Lee, and Icksang Han. In defence of metric learning  
414 for speaker recognition. In *Interspeech*, 2020.
- 415 [12] Alexandre Défossez, Gabriel Synnaeve, and Yossi Adi. Real time speech enhancement in the  
416 waveform domain. In *Interspeech*, 2020.
- 417 [13] Jesse Engel, Lamtharn Hantrakul, Chenjie Gu, and Adam Roberts. Ddsp: Differentiable digital  
418 signal processing. In *International Conference on Learning Representations (ICLR)*, 2020.
- 419 [14] Lei Fan, Qing-Yuan Jiang, Ya-Qi Yu, and Wu-Jun Li. Deep hashing for speaker identification  
420 and retrieval. In *Interspeech*, 2019.
- 421 [15] Electronic Frontier Foundation. Upstream vs. prism. <https://www.eff.org/702-spying>,  
422 2018.
- 423 [16] Barton Gellman and Ashkan Soltani. Nsa surveillance program reaches ‘into the past’ to retrieve,  
424 replay phone calls. *The Washington Post*, 2014.
- 425 [17] Hannah Giorgis. When the fbi spied on mlk. *The Atlantic*, 2021.
- 426 [18] Timothy J Hazen, Wade Shen, and Christopher White. Query-by-example spoken term detection  
427 using phonetic posteriorgram templates. In *IEEE Workshop on Automatic Speech Recognition*  
428 *& Understanding*, 2009.
- 429 [19] Hee Soo Heo, Bong-Jin Lee, Jaesung Huh, and Joon Son Chung. Clova baseline system for the  
430 voxceleb speaker recognition challenge 2020. *arXiv preprint arXiv:2009.14153*, 2020.
- 431 [20] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*,  
432 9(8):1735–1780, 11 1997.

- 433 [21] J. Hsu. pyworld. [https://github.com/JeremyCCHsu/](https://github.com/JeremyCCHsu/Python-Wrapper-for-World-Vocoder)  
434 Python-Wrapper-for-World-Vocoder, 2021.
- 435 [22] Umar Iqbal, Pouneh Nikkhah Bahrami, Rahmadi Trimananda, Hao Cui, Alexander Gamero-  
436 Garrido, Daniel Dubois, David Choffnes, Athina Markopoulou, Franziska Roesner, and Zubair  
437 Shafiq. Your echos are heard: Tracking, profiling, and ad targeting in the amazon smart speaker  
438 ecosystem. *arXiv preprint arXiv:2204.10920*, 2022.
- 439 [23] Qin Jin, Arthur R. Toth, Tanja Schultz, and Alan W Black. Voice convergin: Speaker de-  
440 identification by voice transformation. In *International Conference on Acoustics, Speech, and*  
441 *Signal Processing (ICASSP)*, Apr. 2009.
- 442 [24] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. In  
443 *IEEE Transactions on Big Data*, volume 7, pages 535–547, 2019.
- 444 [25] Takuhiro Kaneko, Hirokazu Kameoka, Kou Tanaka, and Nobukatsu Hojo. Cyclegan-vc3:  
445 Examining and improving cyclegan-vcs for mel-spectrogram conversion. In *Interspeech*, 2020.
- 446 [26] Andre Kassis and Urs Hengartner. Practical attacks on voice spoofing countermeasures. *arXiv*  
447 *preprint arXiv:2107.14642*, 2021.
- 448 [27] Dia Kayyali. The history of surveillance and the black community. [https://www.eff.org/](https://www.eff.org/deeplinks/2014/02/history-surveillance-and-black-community)  
449 [deeplinks/2014/02/history-surveillance-and-black-community](https://www.eff.org/deeplinks/2014/02/history-surveillance-and-black-community), 2014.
- 450 [28] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
451 *arXiv:1412.6980*, 2014.
- 452 [29] Lantian Li, Chao Xing, Dong Wang, Kaimin Yu, and Thomas Fang Zheng. Binary speaker  
453 embedding. In *2016 10th International Symposium on Chinese Spoken Language Processing*  
454 *(ISCSLP)*, 2016.
- 455 [30] Yinghao Aaron Li, Ali Zare, and Nima Mesgarani. Starganv2-vc: A diverse, unsupervised,  
456 non-parallel framework for natural-sounding voice conversion. In *Interspeech*, 2021.
- 457 [31] Zhuohang Li, Yi Wu, Jian Liu, Yingying Chen, and Bo Yuan. Advpulse: Universal,  
458 synchronization-free, and targeted audio adversarial attacks via subsecond perturbations. In  
459 *CCS*, 2020.
- 460 [32] Loren Lugosch, Mirco Ravanelli, Patrick Ignoto, Vikrant Singh Tomar, and Yoshua Bengio.  
461 Speech model pre-training for end-to-end spoken language understanding. In *Interspeech*, 2019.
- 462 [33] Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger.  
463 Montreal forced aligner: trainable text-speech alignment using kaldi. In *Interspeech*, 2017.
- 464 [34] RG McCurdy. Tentative standards for sound level meters. *Electrical Engineering*, 55(3):260–  
465 263, 1936.
- 466 [35] Masanori Morise, Hideki Kawahara, and Haruhiro Katayose. Fast and reliable f0 estimation  
467 method based on the period extraction of vocal fold vibration of singing voice and speech.  
468 *Journal of the Audio Engineering Society*, February 2009.
- 469 [36] Masanori Morise, Fumiya Yokomori, and Kenji Ozawa. World: A vocoder-based high-quality  
470 speech synthesis system for real-time applications. *IEICE Transactions on Information and*  
471 *Systems*, E99.D(7):1877–1884, 2016.
- 472 [37] Max Morrison, Brian Tang, Gefei Tan, and Bryan Pardo. Reproducible subjective evaluation.  
473 In *ICLR Workshop on ML Evaluation Standards*, April 2022.
- 474 [38] A. Nagrani, J. S. Chung, and A. Zisserman. Voxceleb: a large-scale speaker identification  
475 dataset. In *Interspeech*, 2017.
- 476 [39] Krishna Kanth Nakka and Mathieu Salzmann. Learning transferable adversarial perturbations.  
477 In *Neural Information Processing Systems (NeurIPS)*, 2021.

- 478 [40] Patrick O'Reilly, Pranjal Awasthi, Aravindan Vijayaraghavan, and Bryan Pardo. Effective  
479 and inconspicuous over-the-air adversarial examples with adaptive filtering. In *International*  
480 *Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022.
- 481 [41] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: An asr  
482 corpus based on public domain audio books. In *International Conference on Acoustics, Speech,*  
483 *and Signal Processing (ICASSP)*, 2015.
- 484 [42] Manuel Pariente. pystoi. <https://github.com/mpariente/pystoi>, 2021.
- 485 [43] Tae Jin Park, Naoyuki Kanda, Dimitrios Dimitriadis, Kyu J. Han, Shinji Watanabe, and Shrikanth  
486 Narayanan. A review of speaker diarization: Recent advances with deep learning. *Computer*  
487 *Speech Language*, 72, November 2021.
- 488 [44] Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron Courville. Film:  
489 Visual reasoning with a general conditioning layer. In *AAAI*, 2018.
- 490 [45] Jordi Pons, Santiago Pascual, Giulio Cengarle, and Joan Serrà. Upsampling artifacts in neural  
491 audio synthesis. In *International Conference on Acoustics, Speech and Signal Processing*  
492 *(ICASSP)*, 2021.
- 493 [46] A.W. Rix, J.G. Beerends, M.P. Hollier, and A.P. Hekstra. Perceptual evaluation of speech  
494 quality (pesq)-a new method for speech quality assessment of telephone networks and codecs.  
495 In *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2001.
- 496 [47] Damien Ronsson and Milos Cernak. Ac-vc: Non-parallel low latency phonetic posteriorgrams  
497 based voice conversion. In *IEEE Automatic Speech Recognition and Understanding Workshop*  
498 *(ASRU)*, 2021.
- 499 [48] Takaaki Saeki, Yuki Saito, Shinnosuke Takamichi, , and Hiroshi Saruwatari. Real-time,  
500 full-band, online dnn-based voice conversion system using a single cpu. In *Interspeech*, 2020.
- 501 [49] L. Schmidt, M. Sharifi, and I. Lopez-Moreno. Large-scale speaker identification. In *International*  
502 *Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2014.
- 503 [50] Ali Shahin Shamsabadi, Francisco Sepúlveda Teixeira, Alberto Abad, Bhiksha Raj, Andrea  
504 Cavallaro, and Isabel Trancoso. Foolhd: Fooling speaker identification by highly impercep-  
505 tible adversarial disturbances. In *International Conference on Acoustics, Speech and Signal*  
506 *Processing (ICASSP)*, 2021.
- 507 [51] Daniel Stoller, Sebastian Ewert, and Simon Dixon. Wave-u-net: A multi-scale neural network  
508 for end-to-end audio source separation. In *International Society for Music Information Retrieval*  
509 *(ISMIR)*, 2018.
- 510 [52] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Good-  
511 fellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference*  
512 *on Learning Representations (ICLR)*, 2014.
- 513 [53] Cees H. Taal, Richard C. Hendriks, Richard Heusdens, and Jesper Jensen. An algorithm for  
514 intelligibility prediction of time–frequency weighted noisy speech. *IEEE Transactions on Audio,*  
515 *Speech, and Language Processing*, 19(7):2125–2136, 2011.
- 516 [54] Natalia Tomashenko, Xin Wang, Emmanuel Vincent, Jose Patino, Brij Mohan Lal Srivastava,  
517 Paul-Gauthier Noé, Andreas Nautsch, Nicholas Evans, Junichi Yamagishi, Benjamin O'Brien,  
518 Anaïs Chanclu, Jean-François Bonastre, Massimiliano Todisco, and Mohamed Maouche. The  
519 VoicePrivacy 2020 challenge: Results and findings. *Computer Speech & Language*,  
520 74:101362, 2022.
- 521 [55] Chongshun Wang, Dani Yogatama, Adam Coates, Tony Han, Awni Y. Hannun, and Bo Xiao.  
522 Lookahead convolution layer for unidirectional recurrent neural networks. In *ICLR Workshop*,  
523 2016.
- 524 [56] Miao Wang, Christoph Boeddeker, Rafael G. Dantas, and Ananda Seelan. pesq. <https://github.com/ludlows/python-pesq>, 2022.  
525

- 526 [57] Human Rights Watch. Q & a: Us warrantless surveillance under section 702 of the  
 527 foreign intelligence surveillance act. [https://www.hrw.org/news/2017/09/14/  
 528 q-us-warrantless-surveillance-under-section-702-foreign-intelligence-surveillance,](https://www.hrw.org/news/2017/09/14/q-us-warrantless-surveillance-under-section-702-foreign-intelligence-surveillance)  
 529 2017.
- 530 [58] Human Rights Watch. Secret evidence and the threat of more war-  
 531 rantless surveillance. [https://www.hrw.org/news/2018/01/11/  
 532 secret-evidence-and-threat-more-warrantless-surveillance,](https://www.hrw.org/news/2018/01/11/secret-evidence-and-threat-more-warrantless-surveillance) 2018.
- 533 [59] Yi Xie, Zhuohang Li, Cong Shi, Jian Liu, Yingying Chen, and Bo Yuan. Enabling fast and  
 534 universal audio adversarial attack using generative model. In *AAAI*, 2021.
- 535 [60] Weiyi Zhang, Shuning Zhao, Le Liu<sup>3</sup>, Jianmin Li, Xingliang Cheng, Thomas Fang Zheng, and  
 536 Xiaolin Hu. Attack on practical speaker verification system using universal adversarial pertur-  
 537 bations. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*,  
 538 2021.
- 539 [61] Ge Zhu, Fei Jiang, and Zhiyao Duan. Y-vector: Multiscale waveform encoder for speaker  
 540 embedding. In *Interspeech*, 2021.

## 541 Checklist

542 The checklist follows the references. Please read the checklist guidelines carefully for information on  
 543 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
 544 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
 545 the appropriate section of your paper or providing a brief inline description. For example:

- 546 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 547 • Did you include the license to the code and datasets? **[No]** The code and the data are  
 548 proprietary.
- 549 • Did you include the license to the code and datasets? **[N/A]**

550 Please do not modify the questions and only use the provided macros for your answers. Note that the  
 551 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
 552 block and only keep the Checklist section heading above along with the questions/answers below.

### 553 1. For all authors...

- 554 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
 555 contributions and scope? **[Yes]** We claim to have designed a system to make highly  
 556 effective, inconspicuous perturbations that can transfer to unforeseen architectures in  
 557 real-time. 1. We show this in our objective 4.4 and subjective evaluations 4.5
- 558 (b) Did you describe the limitations of your work? **[Yes]** See Section 5
- 559 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See  
 560 Section 5
- 561 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 562 them? **[Yes]**

### 563 2. If you are including theoretical results...

- 564 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- 565 (b) Did you include complete proofs of all theoretical results? **[N/A]**

### 566 3. If you ran experiments...

- 567 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
 568 perimental results (either in the supplemental material or as a URL)? **[Yes]** It will be  
 569 anonymized in the supplemental materials
- 570 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
 571 were chosen)? **[Yes]**

- 572 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
573 ments multiple times)? [Yes] We provide a distribution of subjective quality ratings  
574 for each attack in in 2. We, however, do not provide error bars over random seeds for  
575 training.
- 576 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
577 of GPUs, internal cluster, or cloud provider)? [Yes] Stated in 4.3
- 578 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 579 (a) If your work uses existing assets, did you cite the creators? [Yes] We cite authors of  
580 LibriSpeech, VoxCeleb, ResNetSE34v2, ReSEval and Y-Vector
- 581 (b) Did you mention the license of the assets? [Yes] We mention the licenses for datasets  
582 and models in 4.2 and for ReSEval in 4.5. We note that the Y-Vector model does not  
583 have a License, to our knowledge.
- 584 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]  
585
- 586 (d) Did you discuss whether and how consent was obtained from people whose data you're  
587 using/curating? [N/A]
- 588 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
589 information or offensive content? [N/A]
- 590 5. If you used crowdsourcing or conducted research with human subjects...
- 591 (a) Did you include the full text of instructions given to participants and screenshots, if  
592 applicable? [Yes] In the Supplementary Materials
- 593 (b) Did you describe any potential participant risks, with links to Institutional Review  
594 Board (IRB) approvals, if applicable? [Yes] We obtained IRB approval. There are no  
595 known risks to the participants in our subjective evaluation. 4.5
- 596 (c) Did you include the estimated hourly wage paid to participants and the total amount  
597 spent on participant compensation? [Yes] In C

## 598 A VoiceBox implementation

### 599 A.1 Architecture

600 We provide additional details of the VoiceBox architecture used in our experiments.

601 **Phoneme encoder:** The phoneme encoder consists of a feedforward network with two hidden  
602 layers of size 256 and ReLU activations, followed by two LSTM layers with hidden size 256. Using  
603 the LibriSpeech "train-clean-100" dataset and aligned phoneme labels provided by Lugosch et al.  
604 [32], we train the phoneme encoder with cross entropy loss for 25 epochs using the Adam optimizer  
605 with learning rate 1e-3. During training, a linear layer is appended to the phoneme encoder to render  
606 predictions over phoneme labels; we discard this layer after training.

607 **Spectrogram encoder:** We pass 64-bin mel-spectrograms through a  $1 \times 1$  convolutional layer and  
608 gated linear unit (GLU) operating along the frequency axis, followed by a two linear layers with leaky  
609 ReLU activations and hidden size 512. The convolution layer upsamples the channel (frequency)  
610 dimension to 1024, from which the GLU obtains a 512-dimensional representation to pass to the  
611 linear layers.

612 **Source speaker conditioning:** For each frame of audio, we concatenate the 256-dimensional  
613 output of the phoneme encoder, the 512-dimensional output of the spectrogram encoder, and the one-  
614 dimensional pitch / periodicity / loudness features and linearly project to obtain a 512-dimensional  
615 vector. A fixed 512-dimensional source speaker embedding obtained from the ResNetSE34V2 model  
616 is passed through a 2 linear layers with leaky ReLU activations and hidden size 512, followed by batch  
617 normalization and a final linear layer to obtain a 1024-dimensional vector holding scales and biases  
618 for a feature-wise affine transformation (FiLM) of the encoder output. The affine transformation is  
619 applied and the resulting vector is passed to the bottleneck.

620 **Bottleneck:** Our bottleneck consists of two LSTM layers with hidden size 512. The encoder output  
 621 is passed through the LSTM layers with a skip connection concatenated at the output. Following the  
 622 application of a lookahead convolutional network, the resulting vector is projected linearly to obtain  
 623 a 128-dimensional vector of filter controls for the current frame.

624 **Decoder:** We perform filtering with 128 bands per frame, and set  $L_\infty$  bound  $\epsilon = 0.5$  to further  
 625 constrain deviations from a "neutral" control configuration of  $\vec{1}$ . We experimented with alternative  
 626 methods of regularizing predicted filter controls, such as  $L_2$  normalization, but found that simple  $L_\infty$   
 627 clipping yielded satisfactory results.

628 **Training:** We train both VoiceBox and the universal baseline attack using the Adam [28] optimizer  
 629 with learning rate  $1e-4$  for 10 epochs over the LibriSpeech "train-clean" subset.

## 630 A.2 Auxiliary loss

631 We briefly summarize the multi-resolution STFT loss of Defosséz et al. [12] used in the training of  
 632 our attacks. The loss is given by

$$\mathcal{L}_{aux}(g, u) = \mathcal{L}_{stft}(u, g(u)) + \|u - g(u)\|_1 \quad (5)$$

633 where  $\mathcal{L}_{stft}$  is given by a sum of magnitude and spectral convergence losses computed over  $M$   
 634 spectrogram resolutions:

$$\mathcal{L}_{stft}(a, b) = \sum_{i=1}^M \left[ \mathcal{L}_{sc}^{(i)}(a, b) + \mathcal{L}_{mag}^{(i)}(a, b) \right] \quad (6)$$

$$\mathcal{L}_{sc}(a, b) = \frac{\| |STFT(a)| - |STFT(b)| \|_F}{\| |STFT(a)| \|_F} \quad (7)$$

$$\mathcal{L}_{mag}(a, b) = \| \log |STFT(a)| - \log |STFT(b)| \|_1 \quad (8)$$

635 Here  $\| \cdot \|_F$  and  $\| \cdot \|_1$  refer to the Frobenius and  $L_1$  norms, respectively. Following Defosséz et al.,  
 636 we use  $M = 3$  spectrogram resolutions with hop sizes of 50, 120, and 240 samples, FFT lengths of  
 637 512, 1024, and 2048, and window lengths of 240, 600, and 1200, respectively.

## 638 A.3 Streamer

639 We implement a version of VoiceBox for processing audio streams in real time. All components are  
 640 modified to support the computation of perturbations on overlapping windows of audio while the  
 641 input and output streams process non-overlapping chunks. Where possible, components are compiled  
 642 to TorchScript to improve performance. Three buffers are added to process a stream of audio into  
 643 overlapping chunks:

- 644 • Last Frame Input Buffer: Contains the last 128 samples of the last input audio to be appended  
 645 to the oncoming chunk, before computing the overlapping windows.
- 646 • Lookahead Buffer: Contains overlapping frames to be used as lookahead.
- 647 • Output Buffer: Contains the last 128 samples output by VoiceBox, to be overlap-added with  
 648 the next processed chunk.

649 Code for our streaming implementation can be found at <https://master.d3hvhbnf7qxjtf.amplifyapp.com/>.

Table 2: Top-1 (T-1) and top-10 (T-10) recognition accuracy of the ResNetSE34v2 system on the de-identification attacks described in Section 4.3 when all query audio is passed through a Demucs [12] speech enhancement preprocessing stage

Approach	ResNetSe34V2		+Demucs	
	T-1 ↓	T-10 ↓	T-1 ↓	T-10 ↓
White noise	0.13	0.40	0.02	0.09
Spectral gating	0.02	0.11	0.02	0.11
Universal	0.14	0.22	0.87	0.97
VoiceBox	0.02	0.10	0.04	0.15
No attack	0.97	0.99	0.96	0.99

## 651 B Additional experiments

### 652 B.1 Attack robustness to preprocessing

653 In real-world settings, user audio may travel through various preprocessing stages before reaching a  
 654 speaker recognition model. To simulate the presence of such stages, we perform attacks using the  
 655 experiment configuration discussed in Section 4.3 but pass query utterances through a pretrained  
 656 Demucs [12] speech enhancement model en route to the ResNetSE34v2 recognition system; results  
 657 are reported in table 2. The speaker recognition system maintains its top-1 and top-10 accuracy on  
 658 clean queries. In contrast to Chiquier et al. [9] we do not incorporate the enhancement model into our  
 659 adversarial optimization, meaning that Demucs essentially functions as an unseen adversarial defense.  
 660 We find that VoiceBox loses almost no effectiveness when queries are passed through Demucs; by  
 661 comparison, the universal attack loses most of its effectiveness, as Demucs is able to separate the  
 662 low-magnitude but noisy perturbation from clean speech. Both the white noise and spectral gating  
 663 attacks retain their effectiveness, presumably because they degrade audio beyond Demucs’ capability  
 664 to provide coherent reconstructions of the original speech. We leave as future work the exploration  
 665 of the robustness of VoiceBox against more sophisticated preprocessing pipelines and adversarial  
 666 defenses.

### 667 B.2 Enrollment of adversarial queries

668 It is possible that adversarially-perturbed query utterances extracted from a VoiceBox user’s audio  
 669 stream may themselves be enrolled by a surveilling speaker recognition system. In such cases, there  
 670 are two possibilities:

- 671 1. VoiceBox successfully de-identifies user speech, and adversarial queries are enrolled as a  
 672 separate speaker profile from any existing clean utterances of the user
- 673 2. VoiceBox fails to de-identify user speech, and adversarial queries are incorporated into an  
 674 existing profile of the user containing clean utterances.

675 We examine both scenarios, again using the VoiceBox attack discussed in Section 4.3. Results are  
 676 presented in table 3. We find that while VoiceBox is highly effective at de-identifying users against  
 677 a profile constructed from clean (unperturbed) utterances, its de-identification performance suffers  
 678 significantly under both of the aforementioned conditions. This suggests that additional work is  
 679 required to ensure that a VoiceBox user remains a "moving target" to surveilling speaker recognition  
 680 systems. One possible solution is to leverage targeted rather than untargeted attacks: by proactively  
 681 "spoofing" specific (random) locations in the embedding space, VoiceBox may hamper the creation  
 682 of a single matching enrolled profile.

## 683 C Listening Study

684 The listening study consists of 20 evaluation tasks, where the participants are shown the following  
 685 text, along with 5 four second long audio examples and 5 corresponding sliders with values from 0 to  
 686 100:

Table 3: For sets of 15 clean and adversarial query utterances, we compare the top-1 (T-1) and top-10 (T-10) accuracy of speaker recognition with the ResNetSE34v2 model under three conditions. **Clean profile:** 20 clean utterances are enrolled per speaker. **Mixed profile:** 10 clean and 10 adversarial (VoiceBox) utterances are enrolled per speaker. **Adversarial profile:** 20 adversarial (VoiceBox) utterances are enrolled per speaker.

Query processing	Clean profile		Mixed profile		Adversarial profile	
	T-1	T-10	T-1	T-10	T-1	T-10
VoiceBox	0.02	0.10	0.51	0.80	0.78	0.93
None	0.97	0.99	0.83	0.95	0.09	0.28

687 Listen to all recordings of a person speaking. Then, move the sliders to **rate the quality of each**  
688 **audio file from 0 (worst) to 100 (best)**. The higher-quality audio files are the ones that are  
689 more natural sounding, or have fewer audio artifacts (e.g., clicks, pops, noise, or otherwise  
690 sound 'unnatural'). **Note** - Each slider cannot be moved until its corresponding audio file has  
691 been listened to in its entirety.

692 It takes approximately 15 minutes to complete all the evaluation tasks. Participants are also asked to  
693 complete a listening test before proceeding to the audio evaluation. The listening test involves the  
694 participant listening to two audio files, and reporting the number of tones they heard. For each audio  
695 file, the participants are given 3 tries to correctly report the number of tones played. If a participant  
696 succeeds in the listening test and completes the evaluation tasks, then they are paid 3.00 USD, or  
697 equivalently, 12.00 USD / hour. If the participant fails the listening test, then they are paid 0.50 USD.