



**HAL**  
open science

## **The VoicePrivacy 2020 Challenge: Results and findings**

Natalia Tomashenko, Xin Wang, Emmanuel Vincent, Jose Patino, Brij Mohan Lal Srivastava, Paul-Gauthier Noé, Andreas Nautsch, Nicholas Evans, Junichi Yamagishi, Benjamin O'Brien, et al.

### ► **To cite this version:**

Natalia Tomashenko, Xin Wang, Emmanuel Vincent, Jose Patino, Brij Mohan Lal Srivastava, et al..  
The VoicePrivacy 2020 Challenge: Results and findings. 2021. hal-03332224v1

**HAL Id: hal-03332224**

**<https://hal.science/hal-03332224v1>**

Preprint submitted on 2 Sep 2021 (v1), last revised 26 Sep 2022 (v4)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# The VoicePrivacy 2020 Challenge: Results and findings

Natalia Tomashenko<sup>a,\*</sup>, Xin Wang<sup>b</sup>, Emmanuel Vincent<sup>c</sup>, Jose Patino<sup>d</sup>, Brij Mohan Lal Srivastava<sup>f</sup>, Paul-Gauthier Noé<sup>a</sup>, Andreas Nautsch<sup>d</sup>, Nicholas Evans<sup>d</sup>, Junichi Yamagishi<sup>b,e</sup>, Benjamin O'Brien<sup>g</sup>, Anaïs Chanclu<sup>a</sup>, Jean-François Bonastre<sup>a</sup>, Massimiliano Todisco<sup>d</sup>, Mohamed Maouche<sup>f</sup>

<sup>a</sup>*LIA, University of Avignon, Avignon, France*

<sup>b</sup>*National Institute of Informatics (NII), Tokyo, Japan*

<sup>c</sup>*Université de Lorraine, CNRS, Inria, LORIA, France*

<sup>d</sup>*Audio Security and Privacy Group, EURECOM, France*

<sup>e</sup>*University of Edinburgh, UK*

<sup>f</sup>*Inria, France*

<sup>g</sup>*LPL, Aix-Marseille University, France*

---

## Abstract

This paper presents the results and analyses stemming from the first VoicePrivacy 2020 Challenge which focuses on developing anonymization solutions for speech technology. We provide a systematic overview of the challenge design with an analysis of submitted systems and evaluation results. In particular, we describe the voice anonymization task and datasets used for system development and evaluation. Also, we present different attack models and the associated objective and subjective evaluation metrics. We introduce two anonymization baselines and provide a summary description of the anonymization systems developed by the challenge participants. We report objective and subjective evaluation results for baseline and submitted systems. In addition, we present experimental results for alternative privacy metrics and attack models developed as a part of the post-evaluation analysis. Finally, we summarise our insights and observations that will influence the design of the next VoicePrivacy challenge edition and some directions for future voice anonymization research.

*Keywords:* privacy, anonymization, speech synthesis, voice conversion, speaker verification, automatic speech recognition, attack model, metrics, utility

---

## 1. Introduction

Due to the growing demand for privacy preservation in the recent years, privacy-preserving data processing has become an active research area. One reason for this is the European general data protection regulation (GDPR) in the European Union (EU) law and similar regulations in national laws of many

---

\*Corresponding author

Email address: [natalia.tomashenko@univ-avignon.fr](mailto:natalia.tomashenko@univ-avignon.fr) (Natalia Tomashenko)

countries outside the EU concerning the implementation of the data protection principles when treating, transferring or storing personal data.

Although a legal definition of privacy is missing (Nautsch et al., 2019a), speech data contains a lot of personal information that can be disclosed by listening or by automated systems (Nautsch et al., 2019b). This includes, e.g., age, gender, ethnic origin, geographical background, health or emotional state, political orientations, and religious beliefs. Speaker recognition systems can also reveal the speaker’s identity. Therefore, the increased interest in developing the privacy preservation solutions for speech technology is not surprising. This motivated the organization of the VoicePrivacy initiative which was created to foster the development of privacy preservation techniques for speech technology (Tomashenko et al., 2020a). This initiative aims to bring together a new community of researchers, engineers and privacy professionals in order to formulate the tasks of interest, develop evaluation methodologies, and benchmark new solutions through a series of challenges. The first VoicePrivacy challenge<sup>1</sup> was organized as a part of this initiative (Tomashenko et al., 2020a,c).

Existing approaches to privacy preservation for speech can be broadly classified into the following types (Vincent et al., 2021): deletion, encryption, distributed learning, and anonymization. Anonymization refers to the goal of suppressing personally identifiable information in the speech signal, leaving all other attributes intact. Note, that in the legal community, the term “*anonymization*” means that this goal has been achieved. Here, it refers to the task to be addressed, even when the method being evaluated has failed. Approaches to anonymization include noise addition (Hashimoto et al., 2016), speech transformation (Qian et al., 2017; Patino et al., 2020), voice conversion (Jin et al., 2009; Pobar & Ipšić, 2014; Bahmaninezhad et al., 2018; Yoo et al., 2020; Magariños et al., 2017), speech synthesis (Fang et al., 2019; Han et al., 2020a; Srivastava et al., 2020a), adversarial learning (Srivastava et al., 2019), and disentangled representation learning (Aloufi et al., 2020). In comparison to other types of privacy preservation methods, anonymization is more flexible because it can selectively suppress or keep unchanged certain attributes in speech and it can easily be integrated within existing systems. Despite the appeal of anonymization and the urgency to address privacy concerns, a formal definition of anonymization and attacks is missing. Furthermore, the level of anonymization offered by existing solutions is unclear and not meaningful because there are no common datasets, protocols and metrics. For these reasons, the VoicePrivacy 2020 Challenge focuses on the task of voice anonymization.

The paper is structured as follows. The challenge design, including the description of the anonymization task, attack models, datasets, objective and subjective evaluation methodologies with the corresponding privacy and utility metrics, is presented in Section 2. The overview of the baseline and submitted systems is provided in Sections 3. Objective and subjective evaluation results and their comparison and analysis are presented in Section 4. We conclude and

---

<sup>1</sup><https://www.voiceprivacychallenge.org/>

discuss future directions in Section 5.

## 2. Challenge design

In this section, we present an overview of the official challenge setup: anonymization task, corresponding attack models selected for the challenge, data and evaluation methodology. Also we present an additional attack model developed as a part of post-evaluation analysis (Tomashenko et al., 2020b).

### 2.1. Anonymization task and attack models

Privacy preservation is formulated as a game between *users* who publish some data and *attackers* who access this data or data derived from it and wish to infer information about the users (Tomashenko et al., 2020a; Qian et al., 2018; Srivastava et al., 2020b). To protect their privacy, the users publish data that contain as little personal information as possible while allowing one or more downstream goals to be achieved. To infer personal information, the attackers may use additional prior knowledge.

Focusing on speech data, a given privacy preservation scenario is specified by: (i) the nature of the data: waveform, features, etc., (ii) the information seen as personal: speaker identity, traits, spoken contents, etc., (iii) the downstream goal(s): human communication, automated processing, model training, etc., (iv) the data accessed by the attackers: one or more utterances, derived data or model, etc., (v) the attackers’ prior knowledge: previously published data, privacy preservation method applied, etc. Different specifications lead to different privacy preservation methods from the users’ point of view and different attacks from the attackers’ point of view. An example of a privacy preservation scenario, for the case where speaker identity is considered as personal information that should be protected, is illustrated in Figure 1.

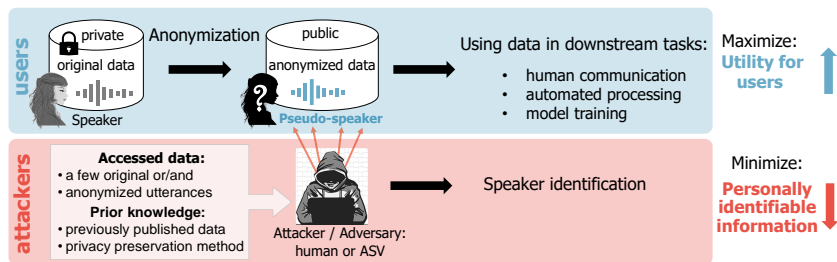


Figure 1: Example of a privacy preservation scenario as a game between *users* and *attackers* for the case where speaker identity is considered as personal information to be protected.

Here, we consider that speakers want to hide their identity while allowing all other downstream goals to be achieved. Attackers want to identify the speakers from one or more utterances.

### 2.1.1. Anonymization task

In order to hide his/her identity, each speaker passes his/her utterances through an anonymization system before publication. The resulting anonymized utterances are called *trial* utterances. They sound as if they were uttered by another speaker, which we call *pseudo-speaker* which may be an artificial voice not corresponding to any real speaker.

The task of challenge participants is to develop this anonymization system. It should: (a) output a speech waveform, (b) hide speaker identity, (c) leave other speech characteristics unchanged, (d) ensure that all trial utterances from a given speaker are uttered by the same pseudo-speaker, while trial utterances from different speakers are uttered by different pseudo-speakers.

The requirement (c) is assessed via a range of *utility* metrics. Specifically, ASR performance using a model trained on original (non-anonymized) data and subjective speech intelligibility and naturalness were measured during the challenge, and additional desired goals including ASR training was assessed in the post-evaluation stage. The requirement (d) is motivated by the fact that, in a multi-party human conversation, each speaker cannot change his/her anonymized voice over time and the anonymized voices of all speakers must be distinguishable from each other. This was assessed via a new, specifically designed metric – *gain of voice distinctiveness*.

### 2.1.2. Attack models

The attackers have access to: (a) one or more anonymized trial utterances, (b) possibly, original or anonymized *enrollment* utterances for each speaker. They do not have access to the anonymization system applied by the user. The protection of personal information is assessed via *privacy* metrics, including objective speaker verifiability and subjective speaker verifiability and linkability. These metrics assume different attack models.

The objective speaker verifiability metrics assume that the attacker has access to a single anonymized trial utterance and several enrollment utterances. Two sets of metrics were computed, corresponding to the two situations when the enrollment utterances are original or they have been anonymized (Section 2.3.1). In the latter case, we assume that the utterances have been anonymized in the same way as the trial data using the same anonymization system, i.e., all enrollment utterances from a given speaker are converted into the same pseudo-speaker, and enrollment utterances from different speakers are converted into different pseudo-speakers. We also assume that the pseudo-speaker corresponding to a given speaker in the enrollment set is different from the pseudo-speaker corresponding to that same speaker in the trial set. In the *post-evaluation* stage, we considered alternative anonymization procedures corresponding to stronger attack models when attackers also have access to anonymized training data and can retrain an automatic speaker verification system using this data. We assume that the training, enrollment and trial data have been anonymized using the same system with different corresponding pseudo-speakers.

For the subjective evaluation (Section 2.3.2), two situations are considered. The speaker verifiability metric assumes that the attacker has access to a single

anonymized trial utterance and a single original enrollment utterance when for the speaker linkability metric, we assume that the attacker has access to several original and anonymized trial utterances.

## 2.2. Datasets

Several publicly available corpora are used for the training, development and evaluation of speaker anonymization systems.

*Training set.* The training set comprises the 2,800 h *VoxCeleb-1,2* speaker verification corpus (Nagrani et al., 2017; Chung et al., 2018) and 600 h subsets of the *LibriSpeech* (Panayotov et al., 2015) and *LibriTTS* (Zen et al., 2019) corpora, which were initially designed for ASR and speech synthesis, respectively. The selected subsets are detailed in Table 1 (top).

Table 1: Number of speakers and utterances in the VoicePrivacy 2020 training, development, and evaluation sets.

Subset		Female	Male	Total	#Utter.	
Training	VoxCeleb-1,2	2,912	4,451	7,363	1,281,762	
	LibriSpeech train-clean-100	125	126	251	28,539	
	LibriSpeech train-other-500	564	602	1,166	148,688	
	LibriTTS train-clean-100	123	124	247	33,236	
	LibriTTS train-other-500	560	600	1,160	205,044	
Development	LibriSpeech dev-clean	Enrollment	15	14	29	343
		Trial	20	20	40	1,978
	VCTK-dev	Enrollment	15	15	30	600
		Trial (common)				695
		Trial (different)				10,677
Evaluation	LibriSpeech test-clean	Enrollment	16	13	29	438
		Trial	20	20	40	1,496
	VCTK-test	Enrollment	15	15	30	600
		Trial (common)				70
		Trial (different)				10,748

*Development set.* The development set involves *LibriSpeech dev-clean* and a subset of the VCTK corpus (Veaux et al., 2019), denoted *VCTK-dev* (see Table 1, middle). With the above attack models in mind, we split them into trial and enrollment subsets. For *LibriSpeech dev-clean*, the speakers in the enrollment set are a subset of those in the trial set. For *VCTK-dev*, we use the same speakers for enrollment and trial and we consider two trial subsets: *common* and *different*. The *common* subset comprises utterances #1 – 24 in the VCTK corpus that are identical for all speakers. This is meant for subjective evaluation of speaker verifiability/linkability in a text-dependent manner. The enrollment and *different* subsets comprises distinct utterances for all speakers.

*Evaluation set.* Similarly, the evaluation set comprises *LibriSpeech test-clean* and a subset of VCTK called *VCTK-test* (see Table 1, bottom).

### 2.3. Utility and privacy metrics

We consider objective and subjective privacy metrics to assess speaker re-identification and linkability. We also propose objective and subjective utility metrics in order to assess the fulfillment of the user goals specified in Section 2.1. Specifically, we consider ASR performance using a model trained on clean data and subjective speech intelligibility and naturalness.

#### 2.3.1. Objective metrics

For objective evaluation of anonymization performance, two systems were trained to assess the following characteristics: (1) speaker verifiability and (2) ability of the anonymization system to preserve linguistic information in the anonymized speech. The first system, denoted  $ASV_{\text{eval}}$ , is an automatic speaker verification (ASV) system, which produces log-likelihood ratio (LLR) scores. The second system, denoted  $ASR_{\text{eval}}$ , is an automatic speech recognition (ASR) system which outputs a word error rate (WER) metric. Both  $ASR_{\text{eval}}$  and  $ASV_{\text{eval}}$  were trained on the *LibriSpeech-train-clean-360* dataset using the Kaldi speech recognition toolkit (Povey et al., 2011). These two models were used in the VoicePrivacy official challenge setup (Tomashenko et al., 2020a).

In addition, for post-evaluation analysis, we trained ASV and ASR systems on anonymized speech data. Both models, denoted  $ASV_{\text{eval}}^{\text{anon}}$  and  $ASR_{\text{eval}}^{\text{anon}}$ , were trained in the same way as  $ASV_{\text{eval}}$  and  $ASR_{\text{eval}}$  respectively<sup>2</sup>.

The  $ASV_{\text{eval}}$  system for speaker verifiability evaluation relies on x-vector speaker embeddings and probabilistic linear discriminant analysis (PLDA) (Snyder et al., 2018). To assess objective speaker verifiability, three metrics are computed in the official challenge setup: the *equal error rate (EER)* and the *log-likelihood ratio (LLR)* based costs  $C_{\text{llr}}$  and  $C_{\text{llr}}^{\text{min}}$ . In addition, we use recently proposed metrics based on *voice similarity matrices* to assess other aspects of anonymization: *de-identification* and *voice distinctiveness*.

As seen in Figure 2, these metrics are computed for 4 evaluation scenarios corresponding to different types of attacks depending on the amount of the attackers’ knowledge. Following terminology of Srivastava et al. (2021), we consider the following conditions:

1. *Unprotected*: no anonymization is performed by the users; attackers have access to original trial and enrollment data.
2. *Ignorant attacker (oa)*: anonymized trial and original enrollment data are used for evaluation. Users anonymize their trial data, but the attackers are unaware about the applied anonymization algorithm and use original data for enrollment.
3. *Lazy-informed (aa)* anonymized trial and enrollment data are used for evaluation. Attackers know what anonymization algorithm is applied by

---

<sup>2</sup>Scripts for training of  $ASR_{\text{eval}}$  and  $ASV_{\text{eval}}$  and for evaluation are provided: <https://github.com/Voice-Privacy-Challenge/Voice-Privacy-Challenge-2020>

users, but do not know the particular parameters, so they anonymized the enrollment data accordingly with the same anonymization system but different parameters.

4. *Semi-informed (aa with the model retrained on anonymized data)*: attackers have the same knowledge as in the previous case (anonymization algorithms, but not the parameters used by given speakers), but in addition to this, attackers also generate anonymized training data and re-train an  $ASV_{eval}$  model on this data. These attackers are supposed to be the strongest ones among the considered in this paper. This evaluation scenario was investigated during the post-evaluation stage.

The number of target and impostor trials in the development and evaluation datasets is given in Table 2.

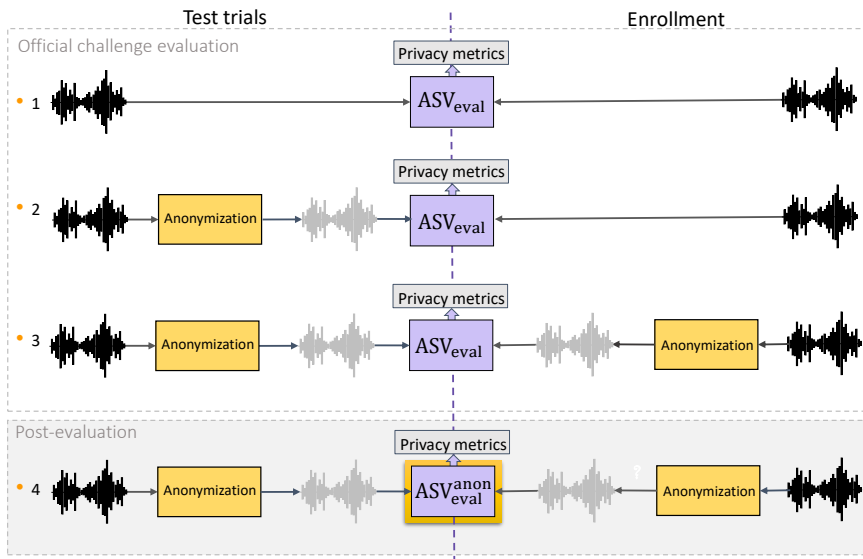


Figure 2: ASV evaluation for the official challenge setup using  $ASV_{eval}$  trained on original data is performed for three cases: (1) *Unprotected*: original trial and enrollment data; (2) *Ignorant attacker (oa)*: anonymized trial data and original enrollment data; (3) *Lazy-informed attacker (aa)*: anonymized trial and enrollment data. ASV evaluation for the post-evaluation analysis is performed using  $ASV_{eval}^{anon}$  trained on anonymized data for case (4) *Semi-informed attacker (aa)*: anonymized trial and enrollment data.

The objective evaluation metrics for privacy and utility are listed below.

**Equal error rate (EER).** Denoting by  $P_{fa}(\theta)$  and  $P_{miss}(\theta)$  the false alarm and miss rates at threshold  $\theta$ , the EER corresponds to the threshold  $\theta_{EER}$  at which the two detection error rates are equal, i.e.,  $EER = P_{fa}(\theta_{EER}) = P_{miss}(\theta_{EER})$ .

**Log-likelihood-ratio cost function ( $C_{llr}$  and  $C_{llr}^{min}$ ).**  $C_{llr}$  is computed from PLDA scores as defined by Brümmer & Du Preez (2006) and Ramos & Gonzalez-



Table 2: Number of speaker verification trials.

Subset		Trials	Female	Male	Total
Development	LibriSpeech dev-clean	Target	704	644	1,348
		Impostor	14,566	12,796	27,362
	VCTK-dev	Target (common)	344	351	695
		Target (different)	1,781	2,015	3,796
		Impostor (common)	4,810	4,911	9,721
		Impostor (different)	13,219	12,985	26,204
Evaluation	LibriSpeech test-clean	Target	548	449	997
		Impostor	11,196	9,457	20,653
	VCTK-test	Target (common)	346	354	700
		Target (different)	1,944	1,742	3,686
		Impostor (common)	4,838	4,952	9,790
		Impostor (different)	13,056	13,258	26,314

Rodriguez (2008). It can be decomposed into a discrimination loss ( $C_{\text{llr}}^{\text{min}}$ ) and a calibration loss ( $C_{\text{llr}} - C_{\text{llr}}^{\text{min}}$ ).  $C_{\text{llr}}^{\text{min}}$  is estimated by optimal calibration using monotonic transformation of the scores to their empirical LLR values.

**Voice similarity matrices.** To visualize anonymization performance across different speakers in a dataset, voice similarity matrices have been proposed by Noé et al. (2020). A voice similarity matrix  $M = (S(i, j))_{1 \leq i \leq N, 1 \leq j \leq N}$  is defined for a set of  $N$  speakers using similarity values  $S(i, j)$  computed for speakers  $i$  and  $j$  as follows:

$$S(i, j) = \text{sigmoid} \left( \frac{1}{n_i n_j} \sum_{\substack{1 \leq k \leq n_i \text{ and } 1 \leq l \leq n_j \\ k \neq l \text{ if } n_i = n_j}} LLR(x_k^{(i)}, x_l^{(j)}) \right) \quad (1)$$

where  $LLR(x_k^{(i)}, x_l^{(j)})$  is the log-likelihood-ratio from the comparison of the  $k$ -th segment from the  $i$ -th speaker and the  $l$ -th segment from the  $j$ -th speaker,  $n_i, n_j$  are the numbers of segments for the corresponding speakers. Three types of matrices are computed:  $M_{oo}$  – on original data;  $M_{aa}$  – on anonymized data; and  $M_{oa}$  – on original and anonymized data. In the latter case, for computing  $S(i, j)$ , we use original data for speaker  $i$  and anonymized data for speaker  $j$ .

**De-identification and voice distinctiveness preservation.** Using voice similarity matrices, two alternative metrics can be computed as was proposed by Noé et al. (2020): de-identification (DeID) and gain of voice distinctiveness ( $G_{\text{VD}}$ ). They are estimated based on the ratio of diagonal dominance for a pair of matrices  $\{M_{oa}, M_{oo}\}$  and  $\{M_{oo}, M_{oo}\}$  correspondingly. The diagonal dominance  $D_{\text{diag}}(M)$  is defined as the absolute difference between the mean values of diagonal and off-diagonal elements:

$$D_{\text{diag}}(M) = \left| \sum_{1 \leq i \leq N} \frac{S(i, i)}{N} - \sum_{\substack{1 \leq j \leq N \\ \text{and } 1 \leq k \leq N \\ j \neq k}} \frac{S(j, k)}{N(N-1)} \right| \quad (2)$$

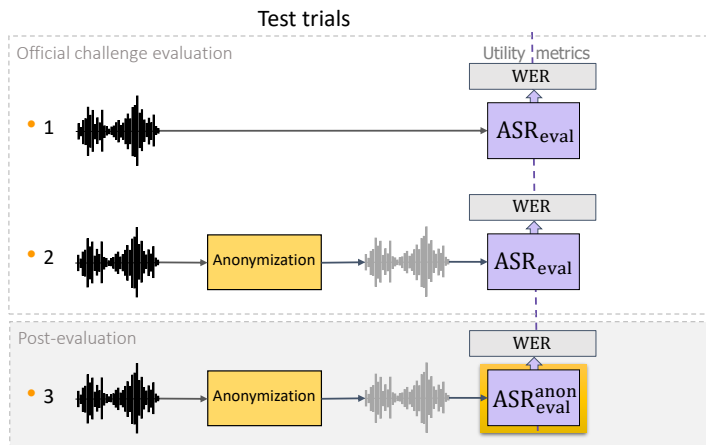


Figure 3: ASR evaluation for the official challenge setup using  $ASR_{eval}$  trained on original data is performed for two cases: (1) original trial data and (2) anonymized trial data. ASR evaluation for the post-evaluation analysis is performed using  $ASR_{eval}^{anon}$  trained on anonymized data for case (3) anonymized trial data.

The *de-identification* is calculated as  $DeID = 1 - D_{diag}(M_{oa})/D_{diag}(M_{oo})$ .  $DeID = 100\%$  assumes perfect de-identification, while  $DeID = 0$  corresponds to a system which achieves no de-identification. *Gain of voice distinctiveness* is defined as  $G_{VD} = 10 \log_{10} (D_{diag}(M_{aa})/D_{diag}(M_{oo}))$ , where 0 means that the voice distinctiveness remains globally the same in the protected space, and gain above or below 0 corresponds respectively to increase or loss of global voice distinctiveness.

**Word error rate (WER).** ASR performance is assessed using  $ASR_{eval}$  which is based on the adapted Kaldi recipe for LibriSpeech involving an acoustic model with a factorized time delay neural network (TDNN-F) architecture (Povey et al., 2018; Peddinti et al., 2015), trained on *LibriSpeech-train-clean-360* dataset, and a trigram language model. As shown in Figure 3, the (1) original and (2) anonymized trial data is decoded using the provided pretrained system  $ASR_{eval}$  and the WERs are calculated. For the post-evaluation analysis, we also perform decoding of anonymized trial data using the  $ASR_{eval}^{anon}$  model trained on anonymized data (Figure 3, case 3).

### 2.3.2. Subjective metrics

We consider two subjective privacy metrics (*speaker verifiability* and *speaker linkability*), and two utility subjective metrics (*speech intelligibility* and *speech naturalness*).

**Subjective speaker verifiability, speech intelligibility, and naturalness.** These three metrics are evaluated in a unified subjective evaluation test illustrated in Figure 4. The input speech trial can be an original or anonymized test

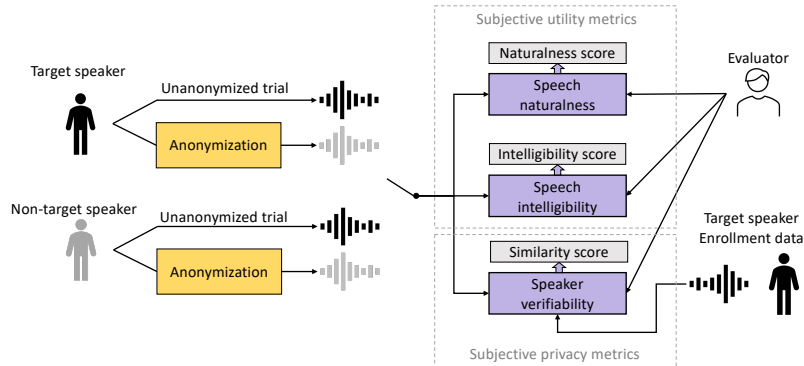


Figure 4: Design of subjective evaluation on speaker verifiability, speech intelligibility, and speech naturalness.

set trial from a target or a non-target speaker. For intelligibility of the input trial, the evaluator assigns a score from 1 (‘totally unintelligible’) to 10 (‘totally intelligible’). For naturalness, the evaluator assigned a score from 1 (‘totally unnatural’) to 10 (‘totally natural’). For verifiability, the evaluator is required to listen to one original enrollment utterance of the target speaker and rate the similarity between the input trial and the enrollment voice using a scale of 1 to 10, where 1 denotes ‘different speakers’ and 10 denotes ‘the same speaker’ with highest confidence. The evaluator were instructed to assign the scores through a role-playing game<sup>3</sup>.

Every evaluator was required to evaluate 36 trials in one session, following the procedures in Figure 4. He or she can also evaluate more than one session. The input trials were randomly sampled from the speakers in the three test sets. The ratio of anonymized and original trials is roughly 1:1, so does the ratio of trials from target and non-target speakers. Among anonymized trials, the ratio of trials from each submitted anonymization system is also balanced. There are 47 native English speakers participated in the evaluation and evaluated 16,200 trials. The decomposed number of trials over the three test sets are listed in Table 3.

***Perception of speaker identity and speaker linkability.*** To evaluate the perception of speaker identity by humans is not simple. In the subjective metrics described previously, we staid as close as possible to the objective metrics. But some questions remain open about potential biases like the memorisation bias (the listener listens to a voice before the current trial) or the priming effect, well known in cognitive psychology. In order both to assess speaker linkability (i.e., the ability to cluster several utterances into speakers) and to decrease as much as possible the influence of perceptual biases, we designed a clustering-based perceptual experiment and the corresponding metrics. A specific software was

<sup>3</sup>Details can be found in (Tomashenko et al., 2021), Section 4.1

Table 3: Number of trials evaluated in subjective evaluation on verifiability, intelligibility, and naturalness. Note that for subjective evaluation, two *LibriSpeech* male speakers are re-used for subjective evaluation. Anonymized trials for subjective evaluation are from 9 anonymized systems (baseline and primary participants’ systems). The number of speakers is 30 (15 male and 15 female) in each dataset.

Test set	Trials	Female	Male	Total
LibriSpeech test-clean	Original	1,330	1,330	2,660
	Anonymized	1,330	1,330	2,660
VCTK-test (common)	Original	1,380	1,380	2,760
	Anonymized	1,380	1,380	2,760
VCTK-test (different)	Original	1,340	1,340	2,680
	Anonymized	1,340	1,340	2,680

developed for this purpose (O’Brien et al., 2021).<sup>4</sup> As perceptual experiments are very demanding in terms of human efforts, we evaluated during this first step only the effects of the two baseline anonymization systems. 74 evaluators were recruited, 29 are native-English speakers and the others are either bilingual or held a high-level of English. Each evaluator did only one session composed of three trials, which gives a great total of 222 trials. Each trial includes 16 different recordings divided between 3 reference speakers and 1 distractor. Adding a distractor helps to verify that the listeners are focused on speaker specificities and are not disturbed by other acoustic differences. The anonymized distractor speaker was used to examine whether anonymization systems affected speaker discrimination performance, e.g. the evaluator either correctly identified the speaker as unique or incorrectly included it in a reference cluster. For a trial, listeners are asked to place a set of recordings from different speakers into 1 to maximum 4 clusters, where clusters represent different speakers, according to subjective speaker voice similarity. In order to decrease as much as possible the potential perceptual biases, during a given session, a speaker is encountered in only 1 trial, and all speakers are of the same gender. Reference speakers are allocated from 2 up to 6 recordings, and the distractor – only 1 recording. For the control trial, genuine speech is used, for other trials half of the recordings are anonymized using the same anonymization system. The data used in the speaker clustering task come from the *VCTK-test (common)* corpus. Unlike all other experiments, only 3 first seconds of each speech recording were used.

As a primary metric, we use macro-average *F-measure* ( $F_1$ ), a classical metric for such a task. We also defined a secondary metric, denoted as *clustering purity*. *Clustering purity* associates each cluster with a single different speaker of a trial and focuses only on precision, compared to  $F_1$  which allows two clusters to be linked to the same speaker and is a harmonic mean of precision and recall.

---

<sup>4</sup><https://demo-lia.univ-avignon.fr/voiceprivacy/instructions>

Clustering *purity* is defined as:

$$purity(C) = \max_{s \in S} \frac{1}{N} \sum_{c \in C} |c \cap s_c|, \quad (3)$$

where  $C$  is the clustering to evaluate,  $c$  is an individual cluster of  $C$ ;  $S$  is the set off all possible combinations of unique speakers assigned to each cluster;  $s_c$  is the speaker label assigned to the cluster  $c$  in the combination  $s$ ; and  $N$  is the number of speech recordings in the trial.

### 3. Anonymization systems

Described in this section are the two baseline systems provided by the challenge organizers as well as those prepared by challenge participants.

#### 3.1. Baseline systems

Two different anonymization systems were provided as challenge baselines<sup>5</sup>. The *primary* baseline, denoted **B1**, is shown in Figure 5. It is inspired from (Fang et al., 2019) and performs anonymization using x-vectors and neural waveform models. It comprises three steps: (1) x-vector (Snyder et al., 2018), pitch (F0) and bottleneck (BN) feature extraction; (2) x-vector anonymization; (3) speech synthesis (SS) using anonymized x-vectors and original F0 and BN features. Step (1) encodes spoken content using 256-dimensional BN features extracted using a TDNN-F ASR AM trained with *LibriSpeech train-clean-100* and *train-other-500* datasets. Speaker encodings are 512-dimensional x-vectors extracted using a TDNN trained with the *VoxCeleb-1,2* dataset. Both extractors are implemented with the Kaldi toolkit. Step (2) computes an anonymized x-vector for every source x-vector. It is generated by averaging a set of  $N^*$  x-vectors selected at random from a larger set of  $N$  x-vectors, itself composed of the  $N$  farthest x-vectors, according to PLDA distances, generated from the *LibriTTS train-other-500* dataset<sup>6</sup>. Step (3) uses a SS AM to generate Mel-filterbank features from the anonymized x-vector and F0+BN features, and a neural source-filter (NSF) waveform model (Wang & Yamagishi, 2019) to synthesize a speech signal from the anonymized x-vector, F0, and Mel-filterbank features. The SS AM and NSF models are both trained with the *LibriTTS train-clean-100* dataset. Full details are available in (Tomashenko et al., 2020c; Srivastava et al., 2020a).

In contrast to the primary baseline, the *secondary* baseline, denoted **B2**, does not require any training data and is based upon traditional signal processing techniques (Patino et al., 2020). It employs the McAdams coefficient (McAdams, 1984) to achieve anonymization by shifting the pole positions derived from the linear predictive coding (LPC) analysis of speech signals. The process is depicted

<sup>5</sup><https://github.com/Voice-Privacy-Challenge/Voice-Privacy-Challenge-2020>

<sup>6</sup>In the baseline, we use  $N = 200$  and  $N^* = 100$ .

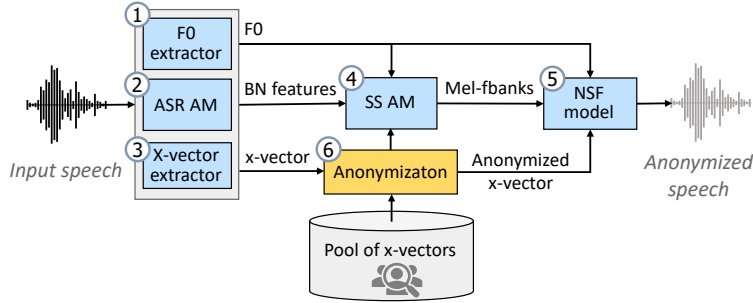


Figure 5: Primary baseline anonymization system (**B1**).

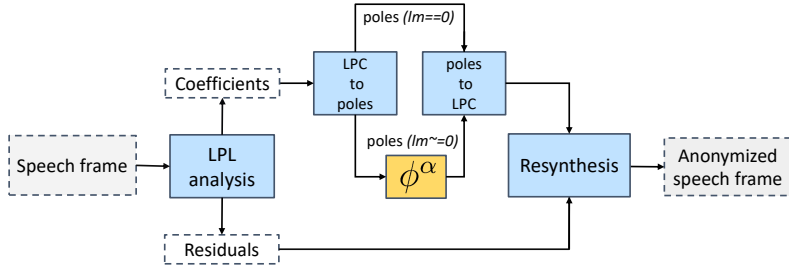


Figure 6: Secondary baseline anonymization system (**B2**).

in Figure 6. It starts with the application of frame-by-frame LPC source-filter analysis to derive LPC coefficients and residuals which are set aside and retained for later resynthesis. The McAdams transformation is then applied to the angles of each pole (with respect to the origin in the  $z$ -plane), each one of which corresponds to a peak in the spectrum (resembling formant positions). While real-valued poles are left unmodified, the angles  $\phi$  of the poles with a non-zero imaginary part (with values between 0 and  $\pi$  radians) are raised to the power of the McAdams coefficient  $\alpha$  so that a transformed pole has new, shifted angle  $\phi^\alpha$ . The value of  $\alpha$  implies a contraction or expansion of the pole positions around  $\phi = 1$ . For a sampling rate of 16kHz, i.e. for data used in the challenge,  $\phi = 1$  corresponds to approximately 2.5kHz which is the approximate mean formant position (Ghorshi et al., 2008). Corresponding complex conjugate poles are similarly shifted in the opposite direction and the new set of poles, including original real-valued poles, are then converted back to LPC coefficients. Finally, LPC coefficients and residuals are used to resynthesise a new speech frame in the time domain. Full details are available in (Patino et al., 2020).

### 3.2. Submitted systems

The VoicePrivacy Challenge attracted 45 participants from both academic and industrial organizations and 13 countries, all represented by 25 teams. Among the 5 allowed submissions by each team, participants were required to designate one as their primary system with any others being designated as

contrastive systems. With full descriptions available elsewhere, we provide only brief descriptions of the 16 successful, eligible submissions, a summary of which is provided in Table 4 which shows system identifiers (referred to below) in column 3. Most systems submitted to the VoicePrivacy 2020 challenge were inspired by the primary baseline (see Section 3.2.1). One submission is based upon the secondary baseline (see Section 3.2.2) whereas two others are not related either (see Section 3.2.3).<sup>7</sup>

Table 4: Challenge submissions, teams names and organizations. Submission identifiers (IDs) for each system are shown in the last column (ID) and comprise: <team id: first letter of the team name><submission deadline<sup>8</sup>: 1 or 2><c, if the system is contrastive><index of the contrastive system>. Blue star symbols  $\star$  in the first column indicate teams submitted the anonymized training data for post-evaluation analysis;  $\textcircled{1}$  and  $\textcircled{2}$  – teams developed their systems from the baseline-1 and baseline-2 respectively, and  $\textcircled{\phantom{0}}$  – other submissions.

Team (Reference)	Organization(s)	Sys.
AIS-lab JAIST (Mawalim et al., 2020) $\textcircled{A}$	<ul style="list-style-type: none"> <li>•Japan Advanced Institute of Science and Technology, Japan</li> <li>•NECTEC, National Science and Technology Development Agency, Thailand</li> </ul>	A1 A2
DA-IICT Speech Group (Gupta et al., 2020) $\textcircled{D}$	<ul style="list-style-type: none"> <li>•Dhirubhai Ambani Institute of Information and Communication Technology, India</li> </ul>	D1
Idiap-NKI (Dubagunta et al., 2020) $\textcircled{I}$	<ul style="list-style-type: none"> <li>•Idiap Research Institute, Martigny, Switzerland</li> <li>•École Polytechnique Fédérale de Lausanne (EPFL), Switzerland</li> <li>•Netherlands Cancer Institute (NKI), Amsterdam, Netherlands</li> </ul>	I1
Kyoto Team (Han et al., 2020b) $\textcircled{K} \star$	<ul style="list-style-type: none"> <li>•Kyoto University, Kyoto, Japan</li> <li>•National Institute of Information and Communications Technology, Kyoto, Japan</li> </ul>	K2
MultiSpeech (Champion et al., 2020b) $\textcircled{M} \star$	<ul style="list-style-type: none"> <li>•Université de Lorraine, CNRS, Inria, LORIA, Nancy, France</li> <li>•Le Mans Université, LIUM, France</li> </ul>	M1 M1c1 M1c2 M1c3 M1c4
Oxford System Security Lab (Turner et al., 2020) $\textcircled{O} \star$	<ul style="list-style-type: none"> <li>•University of Oxford, UK</li> </ul>	O1 O1c1
Sigma Technologies SLU (Espinoza-Cuadros et al., 2020a) $\textcircled{S} \star$	<ul style="list-style-type: none"> <li>•Sigma Technologies S.L.U., Madrid, Spain</li> <li>•Universidad Politecnica de Madrid, Spain</li> </ul>	S1 S1c1 S2 S2c1
PingAn (Huang, 2020)	<ul style="list-style-type: none"> <li>•PAII Inc., Palo Alto, CA, USA</li> </ul>	-

### 3.2.1. Submissions derived from Baseline-1

Teams **A**, **M**, **O** and **S** (see identifiers in column 3 of Table 4 and column 1 of Table 5 submitted systems derived from the primary baseline. Table 5 provides

<sup>7</sup>There is also one non-challenge entry work related to the challenge (Huang, 2020). This team worked on the development of stronger attack models for ASV evaluation.

<sup>8</sup>deadline-1: 8th May 2020; deadline-2: 16th June 2020.

an overview of the modifications made by each team to the baseline modules shown in Figure 5. None of the teams modified the x-vector extraction module (#3 in Table 5), whereas two systems have modifications in the anonymization module (#6) Details of specific modifications are described in the following. We focus first on differences made to specific modules, then on specific system attributes.

Table 5: Summary of the challenge submissions derived from **B1**. ✓ and blue color indicate the components and speaker pool data that were modified w.r.t. **B1**.

Sys.	Description of modifications	1	2	3	4	5	6	Data for speaker pool
		F0	ASR	X-vect.	SS	NSF	Anon.	
A2	using singular value modification						✓	LibriTTS: train-other-500 LibriTTS: train-clean-100
A1	different F0 extractor <sup>9</sup> ; x-vector anonymization using statistical-based ensemble regression modeling	✓					✓	
M1	ASR part to extract BN features for SS models (E2E)		✓		✓	✓		
M1c1	ASR part to extract BN features for SS models (E2E); semi-adversarial training to learn linguistic features while masking speaker information		✓		✓	✓		
M1c2	copy-synthesis (original x-vectors)						✓	
M1c3	x-vectors provided to SS AM are anonymized, x-vectors provided to NSF are original						✓	
M1c4	x-vectors provided to SS AM are original, x-vectors provided to NSF are anonymized						✓	
O1	keeping original distribution of cosine distances between speaker x-vectors; GMM for sampling vectors in a PCA-reduced space with the following reconstruction to the fake x-vectors of the original dimension						✓	LibriTTS: train-other-500
O1c1	<b>O1</b> with forced dissimilarity between original and generated x-vectors						✓	VoxCeleb - 1,2
S1	<b>S1c1</b> applied on the top of the <b>B1</b> x-vector anonymization						✓	
S1c1	domain-adversarial training; autoencoders: using gender, accent, speaker id outputs corresponding to adversarial branches in ANN for x-vector reconstruction						✓	
S2	<b>S2c1</b> applied on the top of the <b>B1</b> x-vector anonymization						✓	
S2c1	<b>S1c1</b> with parameter optimization						✓	

**F0**: Only team **A** (Mawalim et al., 2020) modified the pitch extractor. They replaced the baseline F0 extractor with WORLD (Morise et al., 2016) and by SPTK<sup>10</sup> alternatives. While no significant impact upon ASR performance was observed, SPTK F0 estimation was found to have some impacts, albeit inconsistent, upon the ASV EER. Consequently, the final system used the baseline F0 extractor. Post-evaluation work conducted by other authors (Champion et al.,

<sup>9</sup>Different F0 extractors were used in experiments, but the baseline F0 — in the final **A1**

<sup>10</sup>Speech Signal Processing Toolkit (SPTK): <http://sp-tk.sourceforge.net/>



2020a) showed improved anonymization performance when F0 statistics of the original speaker are replaced with those of a pseudo-speaker, without significant impacts upon ASR performance.

**ASR AM, speech synthesis AM and NSF model.** Instead of the baseline hybrid TDNN-F ASR acoustic model, systems **M1** and **M1c1** (Champion et al., 2020b) used an end-to-end model with a hybrid connectionist temporal classification (CTC) and attention architecture (Watanabe et al., 2017) for BN feature extraction. The SS AM and NFS model were then re-trained using the new BN features. In addition, in the **M1c1** contrastive system used semi-adversarial training to learn linguistic features while masking speaker information.

**X-vector anonymization.** All teams explored different approaches to x-vector anonymization. They are described in the following:

◦**A2.** *Singular value modification* (Mawalim et al., 2020). The singular value decomposition (SVD) of the matrix constructed from the utterance-level speaker x-vectors is used for anonymization. The target x-vector is obtained from the least similar centroid using x-vector clustering. Anonymization is performed through modification of the matrix singular values. A singular value threshold parameter determines the dimensionality reduction used in the modification and determines the percentage of the kept non-zero singular values.

◦**A1.** *Statistical-based decomposition with regression models* (Mawalim et al., 2020). The speaker x-vector is decomposed into high and low-variant components which are separately modified using two different regression models. It is argued that the speaker-specific information is contained more within the low-variant, more stable component, which is hence the component upon which the anonymisation must focus.

◦**O1.** *Distribution-preserving x-vector generation* (Turner et al., 2020). The **B1** baseline performs anonymization through x-vector averaging. As a result, the diversity among anonymized voices is less than that among original voices and observable differences in the distribution of x-vectors between original and anonymized data leaves the anonymization system vulnerable to inversion. The work by Turner et al. (2020) investigated the use of GMMs to sample x-vectors in a PCA-reduced space in a way that retains the original distribution of cosine distances between speaker x-vectors, thereby improving robustness to inversion.

◦**O1c1.** *Forced dissimilarity between original and anonymized x-vectors* (Turner et al., 2020). In a slight variation to the **O1** system, the **O1c1** contrastive system generates a new x-vector in the case that original and anonymized x-vectors are not sufficiently dissimilar.

◦**S1c1** & **S2c1**. *Domain-adversarial training* (Espinoza-Cuadros et al., 2020a). Domain adversarial training is used to generate x-vectors with separate gender, accent, and speaker identification adversarial branches in an autoencoder adversarial network (ANN). For system **S2c1**, the parameters of the adversarial branches are tuned to optimised the trade off between the autoencoder and adversarial objectives.

◦**S1** & **S2**. *Domain-adversarial training on top of B1* (Espinoza-Cuadros et al., 2020a). Primary systems **S1** and **S2** are based upon the application of **S1c1** and **S2c1** contrastive systems anonymized x-vectors generated by the **B1** baseline.

◦**M1c2**. *Copy-synthesis* (Champion et al., 2020b). This contrastive system is essentially the **B1** baseline, but without *explicit* x-vector anonymization, It provides some insights into the added benefit of the latter, beyond simple copy-synthesis.

◦**M1c3**. *Original x-vectors for NSF*. Another contrastive system for which the NSF model receives original x-vectors while the SS AM receives anonymized x-vectors.

◦**M1c4**. *Original x-vectors for SS AM*. A variation on the above contrastive systems whereby the SS AM receives original x-vectors but the NSF model receives anonymised x-vectors.

◦**A** and **O**. *Speaker pool augmentation*. In addition to their respective modifications made to x-vector anonymization, some teams also investigated augmentation of the x-vector pool using additional *LibriTTS-train-clean-100* (team **A**) and *VoxCeleb-1,2* (team **O**) datasets.

### 3.2.2. Submission derived from Baseline-2

◦**D1**. *Modifications to pole radius* (Gupta et al., 2020). Team **D** investigated modifications to the pole radius (distance from the origin) in addition to the shift in phase, as in the original **B2** baseline. This approach distorts formant frequencies with additional changes to the spectral envelope and hence stands to improve anonymization performance. Pole radii are reduced to 0.975 of the original values whereas the McAdam’s coefficient is set to 0.8 as for the **B2** baseline.

### 3.2.3. Other submissions

◦**K2**. *Anonymization using x-vectors, SS models and a voice-indistinguishability metric* (Han et al., 2020b). Similar to the primary baseline **B1**, system **K2** is also based on x-vector anonymization, while the anonymization process and SS models (and corresponding input features) are quite different to those of the **B1** baseline. Other differences include use of the test dataset in creating the speaker pool. The speech synthesis framework uses two modules: (1) an end-to-end AM

implemented with ESPnet<sup>11</sup> which produces a Mel-spectrogram from filterbank features and speaker x-vectors; (2) a waveform vocoder based on the Griffin-Lim algorithm (Griffin & Lim, 1984) which produces a speech waveform from the Mel-spectrogram after conversion to a linear scale spectrogram. A voice indistinguishability metric (Han et al., 2020a) inspired by differential privacy concepts (Dwork, 2009), is applied during x-vector perturbation to select target speaker x-vectors.

◦**I1**. *Modifications to formants, F0 and speaking rate* (Dubagunta et al., 2020). The **I1** system is based upon a signal-processing technique inspired from (van Son, 2020). VTL characteristics are modified by adjusting playback speed to linearly shift formant frequencies. Individual formants are then shifted to specific target values chosen from a set of randomly chosen speakers in the *LibriSpeech-train-other-500* dataset. F0 and the speaking rate are also adjusted using a *pitch synchronous overlap and add method* (Moulines & Charpentier, 1990). Additional processing includes the exchange of F4 and F5 bands using a Hann filter method and the addition of modulated pink noise to the speaker F6–F9 bands for formant masking.

## 4. Results

In this section we report objective and subjective evaluation results for the systems described in Section 3.

### 4.1. Objective evaluation results

Objective evaluation aims to gauge how well each system fulfills the requirements formulated in Section 2.1.1.

#### 4.1.1. Privacy: objective speaker verifiability

Speaker verifiability results are illustrated in Figure 7. Results show EERs averaged across all test datasets for the *ignorant (oa)* and *lazy-informed (aa)* attack models described in Section 2.1.2. Without anonymization, the EER is 3.29%. Anonymization is expected to provoke increases in the EER. When only trial data is anonymized (**oa** condition, light bars in Figure 7), the EER increases for all anonymization systems: from 22.56% for **M1c4** to 53.37% for **M1c1**. Better anonymization is achieved by using x-vector based anonymization systems (**K2**, **A\***, **S\***, **M\***, **B1**, **O\***), than by using pure signal-processing ones (**B2**, **D1**, **I1**). **M1c2** and **M1c4** systems do perform worse, though they provide non-anonymized x-vectors to the speech synthesis AM; they are not expected to perform well. Thus, degradation in EER for these two systems may be caused not by explicit suppression of the speaker-related information (that is assumed to be represented by speaker x-vectors), but by other factors, in particular, by acoustic mismatch between original and synthesised speech.

---

<sup>11</sup><https://github.com/espnet/espnet/tree/master/egs/librispeech/tts1>

Systems **K2**, **A\***, **M1c1**, **M1**, **B1** all produce EERs above 50%, indicating that the anonymization requirement is fully met.

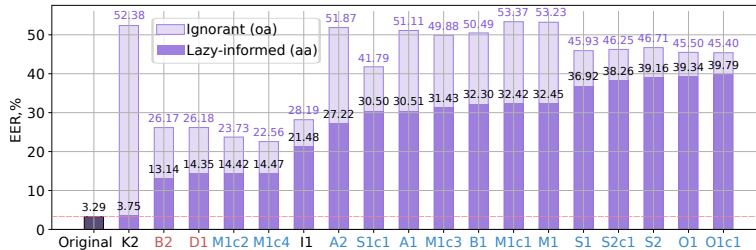


Figure 7: Mean EER values over all test datasets for different anonymization systems and for original data; **oa** denotes the evaluation case #2 in Figure 2 (original enrollment and anonymized trial data); and **aa** – case #3 (both enrolment and trial data are anonymized). Blue and red colors in the system notations indicate systems developed from **B1** and **B2**, respectively. Larger EER values correspond to better privacy preservation.

However, anonymization of both enrollment and trial utterances (Figure 7, **aa** condition, dark bars) demonstrates universally lower EERs for all systems. While system rankings are different for each attack model, the general trend is consistent; **B1** based systems outperform other methods. Some results are of particular interest. The EER of 3.75% for system **K2** is only marginally above the original EER of 3.29%, despite it being the 3rd best performing system for the **oa** condition. Even the best **O1c1** system achieves an EER of only 37.79%, which is far away from the 50% which indicates successful anonymization. These results highlight the importance of designed anonymization systems under the assumption of a strong attack model. Without them, results may provide a false sense of protection.

For other privacy metrics, we observed a consistent correlation in results with EER (see, for example, Figure 8 which illustrates EER vs  $C_{llr}^{\min}$  results for ignorant and lazy-informed attack models for different datasets and systems).

Due to the space constrains, we will focus on EER metric in this paper (see results for other metrics in (Tomashenko et al., 2021), Section 3).

For the **oa** condition, systems **A1**, **A2**, **M1**, **M1c1** all outperform the **B1** baseline, whereas systems **S2**, **S2c1**, **O1**, **O1c1** all outperform the **B1** baseline for the **aa** condition. There is no intersection between the two system groups and no single system works better than others for both conditions. This observation shows the difficulty in designing and optimising an anonymization system that works well under different attack scenarios. Results for the system **K2** are also of note. **K2** has very high anonymization performance for the **oa** scenario yet very poor anonymization performance for the **aa** condition. This may be explained by the strategy used for x-vector anonymization and TTS system. The anonymized utterances are all acoustically very different in comparison to the original ones. Thus, in the **oa** condition, the EER is high. Instead of being generated from a corpus of recordings containing data collected from a large

number of speakers (relative to the evaluation dataset), system **K2** generates anonymized x-vectors from the evaluation dataset itself. This explains why in the **aa** condition, we can observe distinct confusions between some speakers. However, the number of such confusions is very low, especially for some test sets (see, for example, as a complementary illustration a speaker similarity matrix  $M_{aa}$  for female speakers of the *LibriSpeech-test* set in Figure 13h).

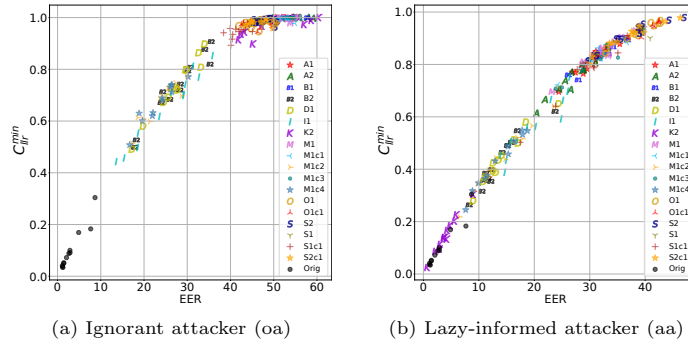


Figure 8: **EER** vs  $C_{lr}^{\min}$  results for two attack models. Each point in the figure corresponds to results on a particular dataset from the set of all 12 VoicePrivacy development and evaluation datasets for a particular system. Higher EER and  $C_{lr}^{\min}$  values correspond to better privacy.

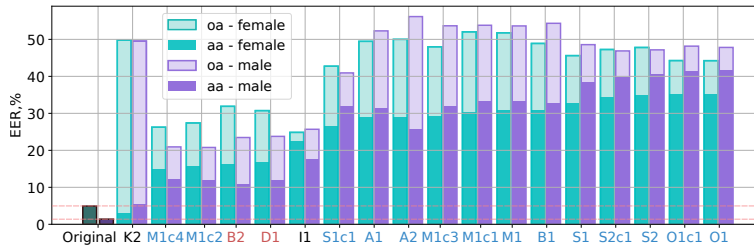


Figure 9: Mean EER values over all datasets for male and female speakers.

Anonymization performance was also found to differ at the gender level. Gender-dependent results for both evaluation conditions averaged across all datasets are illustrated in Figure 9. With no anonymization, the EER is lower for male speakers than for female speakers. With only few exceptions (i.e. **A2**), the opposite is observed after anonymization is applied using systems x-vector based anonymization systems; EERs are generally higher for female speakers than for male speakers. Systems **M1c2** and **M1c4**, for which the SS AM are fed with original x-vectors, are two of the exceptions, indicating that gender-dependent differences are the result of x-vector anonymization rather than any extraneous influence, e.g. acoustic mismatch between original and synthesised data. In contrast, signal-processing approaches show the same gender-dependent trend

as observed for the original data. Nonetheless, depending on the type of the method, the speaker gender seems to influence anonymization performance.

#### 4.1.2. Utility: speech recognition error

Figure 10 shows ASR performance in terms of WER. Since we observed substantial disparities in results, they are illustrated separately for the *LibriSpeech-test* and *VCTK-test* datasets for which WERs are 4.14% and 12.81% respectively for original data (no anonymization). There is a clear explanation for differences in the WER for each dataset; with the ASR system being trained using the *LibriSpeech-train-clean-360* data, performance is better for the matched *LibriSpeech-test* dataset than for the mis-matched *VCTK-test*.

All approaches to anonymization result in increased WERs; any improvement in privacy comes at the expense of utility. It is clear that degradations in the WER are more substantial for the *LibriSpeech-test* dataset than for the *VCTK-test* dataset. The relative increase in WER is 40%–217% for *LibriSpeech-test* and 14–120% for *VCTK-test*.

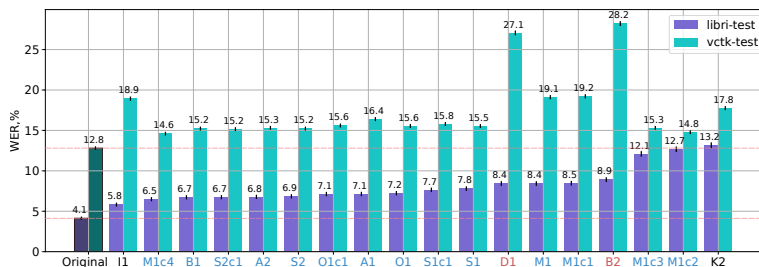


Figure 10: Word error rate (WER) results (with 95% confidence intervals) on *LibriSpeech-test* and *VCTK-test* for different anonymization systems. Smaller values correspond to better utility.

After anonymization, the best WER of 5.83% for the *LibriSpeech* dataset is obtained by the signal-processing based system **I1**. Compared to other systems, however, it performs poorly for the *VCTK-test* dataset. Other signal-processing systems based upon the **B2** baseline fair even worse for this dataset. In average, for both test sets, better results than for other systems (and very close to each other) are obtained using x-vector based anonymization techniques related to the primary baseline: **B1**, **S2c1**, **A2**, **S2**.

Of note is the high WER for the *LibriSpeech-test* dataset and system **M1c2** which performs copy synthesis without x-vector anonymization. Systems **M1c3** and **M1c4**, which also retain original x-vectors, also provoke high WERs for the *LibriSpeech-test* dataset. This finding suggests that resynthesis by itself causes non-negligible degradation to ASR performance. Results for systems **M1** and **M1c1** (vs **B1**), indicate that using the end-to-end acoustic model for BN feature extraction degrades ASV performance for both datasets. Speech recognition accuracy degradation for signal-processing based techniques (**I1**, **D1**, **B2**) is

more stable for different data sets, while for all x-vector based techniques there is a huge gap in performance depending on the datasets, and increase in WER is larger for in-domain data with respect to the data used for training of the  $ASV_{eval}$  model (*LibriSpeech*).

#### 4.1.3. Using anonymized speech data to train attack models

For evaluations reported above, the ASV is trained on regular, non-anonymized data. Reported here are evaluation results for almost identical setup, except that the ASV system is trained using anonymized data according to the semi-informed attacker scenario described in Section 2.3.1. Four teams submitted anonymized *LibriSpeech-train-clean-360* training dataset for primary systems **O1**, **M1**, **S2**, **K2**, and we trained four new corresponding  $ASV_{eval}^{anon}$  models on this data. In addition, we trained two  $ASV_{eval}^{anon}$  models on the training data anonymized by the baseline systems **B1**, **B2**. Models were trained in the same way as before, and have the same topology as the  $ASV_{eval}$  model trained on original data.

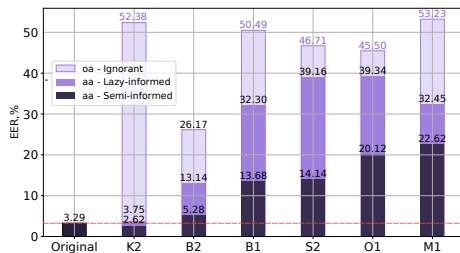


Figure 11: Mean EER,% results over all the test data sets for three types of attack models corresponding to different anonymization systems and for original evaluation without privacy preservation.

Figure 11 shows mean EER results over all VoicePrivacy test data sets for semi-informed attack models (darker, lower bars), lazy-informed, and ignorant attack models. For all anonymization systems and all datasets, training ASV evaluation models on anonymized speech data significantly decreases the EER: EERs are substantially lower for semi-informed than for ignorant and lazy-informed attack models. Thus, the assessment of anonymization systems performed with original data leads to high EERs that give a false impression of protection; if the ASV system is retrained with similarly anonymized data, then ASV performance is much closer to that for original data without any anonymization at all.

#### 4.1.4. Using anonymized speech data for ASR training

Figure 12 shows WERs for ASR systems trained on original data ( $ASR_{eval}$ ) and anonymized speech data ( $ASR_{eval}^{anon}$ ). WERs for  $ASR_{eval}^{anon}$  (lower, darker bars, **a**) are consistently lower than for  $ASR_{eval}$  (upper, lighter bars, **o**). In some cases, WERs decrease to a level close to that of the original system trained on original data (no-anonymization). This finding implies that degradations

to utility can be offset simply by retraining using similarly anonymized data. It improves substantially the trade-off between privacy and utility; there is potential to protect privacy with only modest impacts upon utility.

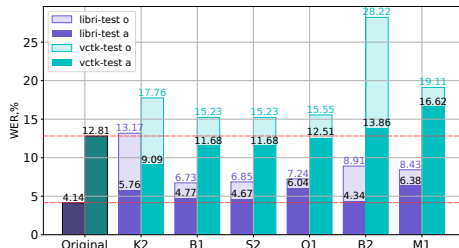


Figure 12: Using anonymized speech data for ASR training on *LibriSpeech-test* and *VCTK-test* for different anonymization systems. **o** – decoding by  $ASR_{eval}$ , **a** – by  $ASR_{eval}^{anon}$ .

#### 4.1.5. Voice distinctiveness preservation

In this section, we consider metrics that are based on voice-similarity matrices, and mainly, *gain of voice distinctiveness* ( $G_{VD}$ ) which evaluates the requirement (d) formulated in Section 2.1.1. In addition, we consider its relation to another metric estimated from similarity matrices – *de-identification* (DeId).

Voice similarity matrices illustrated in Figure 13 show substantial differences between the different approaches. For  $M_{oo}$ , a distinct diagonal in the similarity matrix points out the speaker discrimination ability in the original set, while in  $M_{oa}$ , the diagonal disappears if the protection is good. In  $M_{aa}$ , the diagonal of the matrix emerges, if the resulting pseudo-voices can be distinguished (Noé et al., 2020). The matrices for signal-processing based approaches exhibit a distinct diagonal for  $M_{aa}$  matrices, indicating that voices remain distinguishable after anonymization. Among x-vector based systems, a distinct diagonal for  $M_{aa}$  is observed only for system **K2**.

The scatter plots in Figure 14 show the *gain of voice distinctiveness* ( $G_{VD}$ ) against *de-identification* performance (DeID) for the *LibriSpeech-test* (left) and *VCTK-test* (right) datasets.<sup>12</sup> As described in Section 2.1.1, DeID should be as high as possible, while  $G_{VD}$  should decrease as little possible; the nearer to the top-right of the scatter plots in Figure 14, the better the anonymization.

Results show that systems based upon the **B1** baseline provide close to perfect de-identification, whereas signal-processing based solutions tend to better preserve voice distinctiveness. For the latter, de-identification performance varies between the datasets. Only system **K2** achieves high de-identification with only modest degradation to voice distinctiveness. Results for systems **M1c4** and **M1c2** which use original x-vectors show that copy-synthesis alone also degrades voice distinctiveness. Interestingly, de-identification performance for both systems is comparable to that for signal-processing based methods. These observations are consistent with EER and  $C_{llr}^{min}$  results.

<sup>12</sup>For more details, see Sections 3.4 and 3.5 in (Tomashenko et al., 2021)



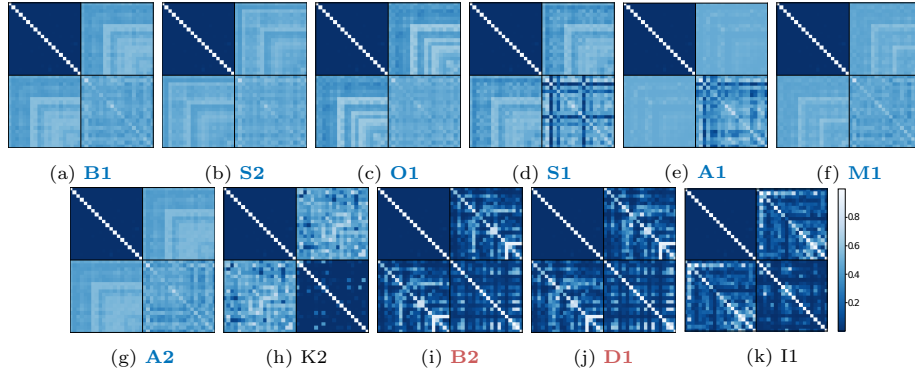


Figure 13: **Voice similarity matrices** for primary systems on *LibriSpeech-test* for female speakers. A global matrix  $M$  for each system is composed of the three matrices  $M_{oo}$ ,  $M_{oa}$  and  $M_{aa}$  as described in Section 2.3.1:  $M = \begin{pmatrix} M_{oo} & M_{oa} \\ M_{oa} & M_{aa} \end{pmatrix}$ .

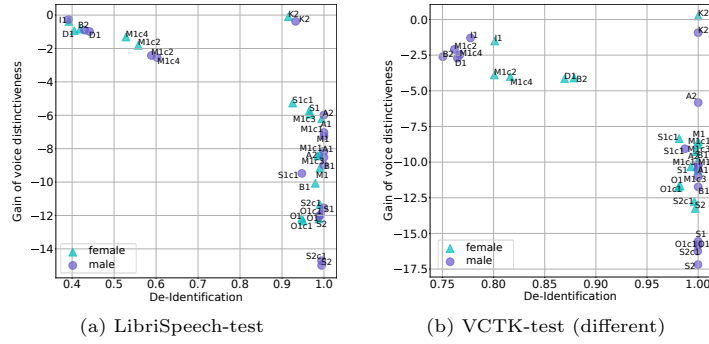


Figure 14: **De-identification (DeID) vs gain of voice distinctiveness ( $G_{VD}$ )**. Higher DeID corresponds to better privacy, higher  $G_{VD}$  – to better distinctiveness of anonymized voices.

Results in Figure 14 also show that different systems lead to differences in voice distinctiveness for different genders. In particular, systems **S2**, **S2c1** preserve distinctiveness better for female speakers than for male speakers, whereas system **A2** better preserves distinctiveness for male speakers.

#### 4.1.6. Relation between privacy and utility metrics

As we observed above, all anonymization systems reduce the utility of speech data. Therefore, it is important to consider the trade-off between privacy and utility. Figure 15 demonstrates the relation between privacy and utility for objective evaluation metrics in the form of scatterplots with (WER, EER-aa), % values for different anonymization systems. The best anonymization system should have maximum EER, and minimum WER (be close to the top-left corner

of Figure 15). We can see that there is no system which provides the best results for both metrics. For *LibriSpeech-test*, best anonymization is achieved using x-vector based anonymization techniques, while the lowest WER corresponds to the system **I1** which is a signal processing approach based on formant shifting. However for *VCTK-test*, results are different for this method, and better results for both metrics were obtained using x-vector-based anonymization techniques.

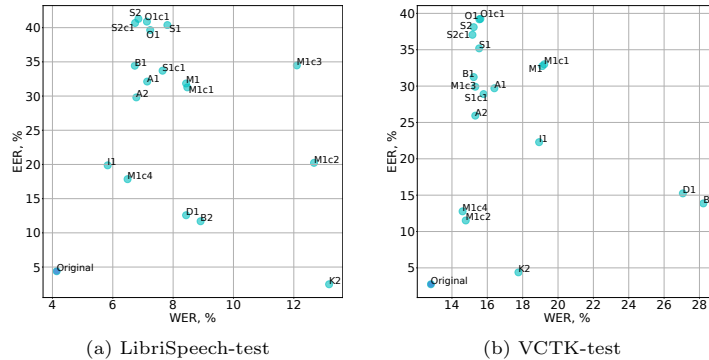


Figure 15: Utility vs privacy: WER,% and EER,% results on anonymized data, EER results correspond to evaluation with lazy-informed attack models. Each point in the figure corresponds to (WER, EER) for a particular anonymization system on a given dataset. Higher EER corresponds to better privacy, lower WER – to better utility.

## 4.2. Subjective evaluation results

This section presents subjective evaluation results for speaker verifiability, speech naturalness, speech intelligibility (Sections 4.2.1–4.2.2), and speaker linkability (Section 4.2.3).

### 4.2.1. Score distribution in violin plot

To reduce perceptual bias of each individual evaluator, naturalness, intelligibility, and verifiability scores from the unified subjective test was processed using normalized-rank normalization (Rosenberg & Ramabhadran, 2017). The processed scores are float numbers varying from 0 to 1. Mann-Whitney-U test was further used for statistical significance tests<sup>13</sup> (Rosenberg & Ramabhadran, 2017).

The score distributions pooled over the three test sets are displayed in Figure 16 as violin plots (Hintze & Nelson, 1998). There are four types of trials: original or anonymized trials from target or non-target speakers. When displaying the results of naturalness and intelligibility, we merge the anonymized trials of both target and non-target speakers. It is only on similarity do we need to separate them so that we can tell how well the anonymization system

<sup>13</sup>Significance test results are reported in (Tomashenko et al., 2021), Tables 16 and 17.

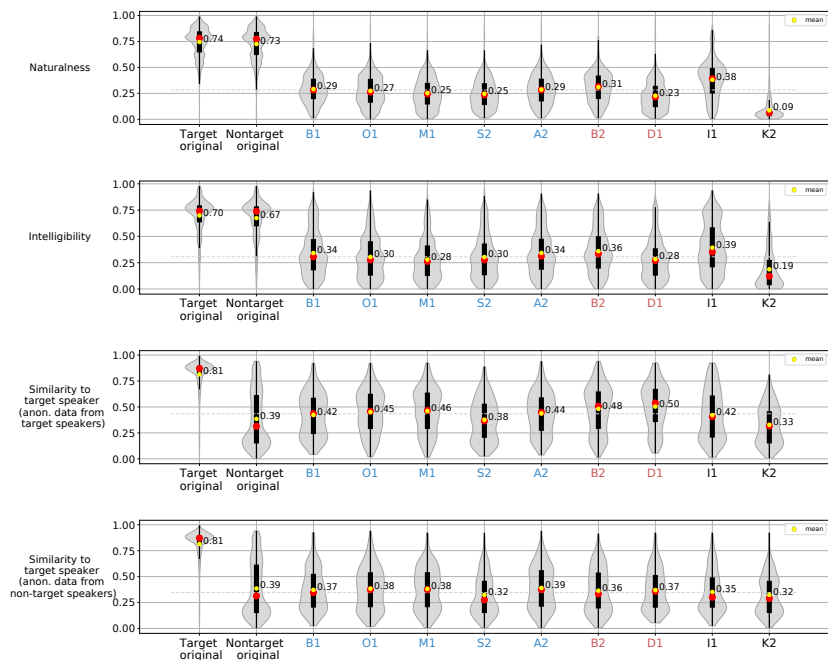


Figure 16: Violin plots of **subjective speech naturalness, intelligibility, and speaker similarity** obtained from the normalized scores. For naturalness and intelligibility, scores from target and non-target anonymized data are pooled; for similarity, scores for anonymized target and non-target speakers data are separately plotted in 3rd and 4th sub-figures, respectively. Dot line indicates median for **B1**. Numbers indicate mean values. Higher values for naturalness and intelligibility correspond to better utility, and lower scores for similarity to target speaker with anonymized data from target speaker – to better privacy.

anonymize the speech of the target speakers. This is the reason why there are four sub-figures in Figure 16. Note that we display the scores of target and non-target original trials separately.

The results of naturalness and intelligibility are as expected. Anonymized samples from all the systems are inferior to the target and non-target original data, and the differences are statistically significant at  $p \ll 0.01$ . This performance gap exists in both methods based on the primary baseline (**B1**, **O1**, **M1**, **S2**, and **A2**) and secondary baseline (**B2**, **D1**). While **I1** outformed other anonymization systems on naturalness, they are still far from perfect on naturalness and intelligibility. More efforts are necessary to address the degradation caused by existing anonymization methods.

On similarity, anonymized data from target speakers are perceptually less similar to the enrollment data of the target speaker than the unanonymized trial of that target speaker. This performance is welcome because it indicates that all the systems achieved a certain degree of anonymization in human perception.

#### 4.2.2. DET curves

To investigate the difference across systems quantitatively, we compute DET curves based on the score distribution. Since there are four types of scores, i.e., {target original, non-target original, target anonymized, non-target anonymized}, we computed the DET curves in the following ways:

- Naturalness and intelligibility DET curves: positive class (anchor) is “target original”, negative class is either “non-target original” or “**anonymized (target and non-target)**” from one anonymization system;
- Similarity DET curve 1: positive class is “target original” or **target anonymized** from one anonymization system, negative class (anchors) is “non-target original”;
- Similarity DET curve 2: positive class is “target original” or **non-target anonymized** from one anonymization system, negative class (anchors) is “non-target original”.

For naturalness and intelligibility, an ideal anonymization system should have a DET curve close to that of original data, indicating similar naturalness and intelligibility scores to the original data and therefore minimum degradation on naturalness and intelligibility. For similarity curve 1, an ideal anonymization system should have a DET curve close to the diagonal line from bottom-right to top-left, indicating that the anonymized data of a target speaker sounds similar to the non-target data.

The four types of DET curves are plotted in Figure 17. As the top two sub-figures demonstrate, the DET curves of the original data are straight lines across the (50%, 50%) point, indicating that the scores of non-target original data are similar to those of the target original data. This is expected because original data should have similar naturalness and intelligibility no matter whether they are from target or non-target speakers. In contrast, the DET curves of anonymized systems are not close to the curve of original data, suggesting that anonymized data are inferior to the original data in terms of naturalness and intelligibility, similar to the messages from the violin plot in previous section.

Among the anonymized systems, the naturalness DET curve of **I1** and **K2** seem to deviate from other systems. While other systems are based on either **B1** or **B2**, **I1** uses a different signal-processing-based approach to change the speech spectra, and **K2** uses a different deep learning method. **I1**’s framework avoids several types of errors such as speech recognition in **B1**, which may contribute to its performance. However, it is interesting to note how different signal processing algorithms result in different perceptual naturalness and intelligibility. Also note that none of the system except **I1** outperformed **B2**.

On similarity, let us focus on the case where target speaker data is anonymized (left-bottom figure in 17). We observe that the DET curve of original data is closer to the bottom-left corner while those of anonymized data are close to top-right corner. In other words, the anonymized data of target speaker produced similar perceptual scores to the non-target speaker data, indicating that

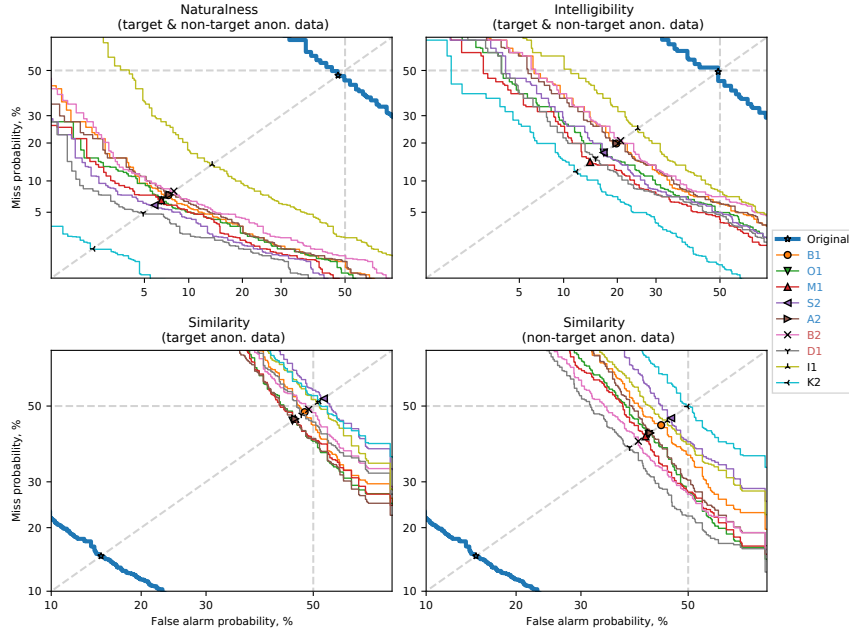


Figure 17: DET curves based on subjective evaluation scores pooled over *LibriSpeech-test* and *VCTK-test* data sets.

anonymized target speaker data sound less similar to the original target speaker. Similar results can be observed from the curves which are separately plotted on the three test sets (see Figure 26 in (Tomashenko et al., 2021)).

The similarity DET curves of **K2**, **S2**, and **I1** on target speaker data seem to be closer to the (50%, 50%) point than others (left-bottom figure in 17). However, the three systems are quite different in terms of naturalness and intelligibility, particularly with **I1** and **K2** achieving the highest and lowest median MOS result, respectively. It implies that an anonymized trial may sound like the voice of a different speaker simply because of the severe distortion caused by anonymization.

In summary, all the submitted the anonymized systems can anonymize the perceived speaker identity to some degrees. However, none of them can produce anonymized trial that is as natural and intelligible as original speech data. One signal-processing-based anonymization method may degrade the naturalness and intelligibility of anonymized trials less severely, but it still introduces degradation on naturalness and intelligibility.

#### 4.2.3. Perception of speaker identity and speaker linkability

We report speaker linkability results for two baseline systems in terms of three metrics: *F-measure* ( $F_1$ ), *classification change* ( $CC$ ), and *clustering purity*. To measure the effects of anonymization speech for each evaluator, we

calculated the *difference* between trial performance on original data and average performance on anonymized trials across all metrics.

We observed a main effect for the mean  $F_1$  difference on a listener native language  $F_{1,64} = 6.5$ ,  $p < 0.05$ ,  $\eta_p^2 = 0.09$ , but no effects on an anonymization system nor speaker gender,  $p > 0.05$ . **B1** evaluators had a greater mean  $F_1$  differences ( $0.24 \pm 0.02$ ) in comparison to **B2** evaluators ( $0.21 \pm 0.02$ ). Post-hoc t-tests showed that non-native English speaking evaluators were more affected by linking natural and anonymized speech recordings ( $0.26 \pm 0.02$ ) in comparison to native English speaking evaluators ( $0.19 \pm 0.022$ ) (Figure 18a).

For the mean CC difference, we found a main effect on speech recording gender  $F_{1,64} = 4.45$ ,  $p < 0.05$ ,  $\eta_p^2 = 0.06$ , and interactions for system  $\times$  language  $F_{1,64} = 4.26$ ,  $p < 0.05$ ,  $\eta_p^2 = 0.06$  and system  $\times$  language  $\times$  gender  $F_{1,64} = 8.75$ ,  $p < 0.01$ ,  $\eta_p^2 = 0.11$ . Post-hoc t-tests revealed that evaluators had a greater mean CC difference when presented male speech recordings ( $0.07 \pm 0.03$ ) in comparison to female ( $-0.03 \pm 0.04$ ) (Figure 18b), as well as native-English speaking evaluators had a greater mean CC difference than **B2** evaluators (Figure 18c). These results suggest that evaluators were able to use the anonymized speech recordings to aid their performance when grouping female speech recordings, whereas performance diminished when they listened to anonymized male speech recordings. Non-native English speaking evaluators lowered their accuracy when presented anonymized stimuli from either system. The different results that we have presented suggest that the effectiveness of an anonymization system can change depending on its users as well as on the evaluators.

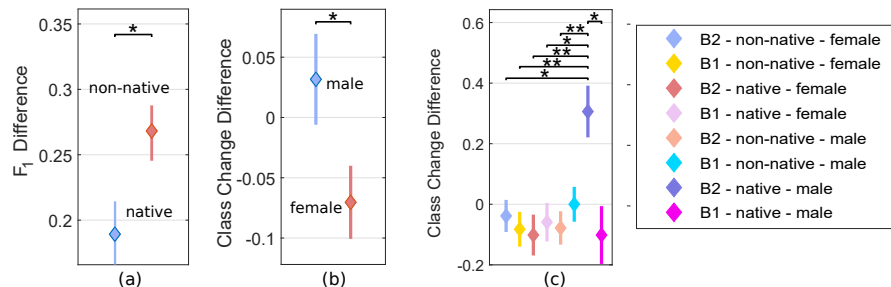


Figure 18: Diamonds and vertical lines represent the means and standard errors, respectively. (a) Mean  $F_1$  difference between native and non-native English speaking evaluators. {\*} signifies  $p < 0.05$ . (b) Mean class change difference between speech recording genders. (c) Mean class change difference between system  $\times$  language  $\times$  gender interactions. {\*, \*\*} signify  $p < \{0.05, 0.01\}$ .

The distribution of clustering purity between different types of trials is displayed in Figure 19a where we can observe higher results for original trials. The Mann-Whitney test shows an effect of the type of trial (original (control) vs anonymized speech (evaluation)) on the clustering purity:  $\chi^2 = 82,688$  ( $p < 0.001$ ) for female speakers and  $\chi^2 = 41,344$  ( $p < 0.001$ ) for male speakers, which indicates different distribution between the original trials and anonymized trials. The distribution of the clustering purity is similar to  $F_1$  for all trial types (see

for example, cumulative distributions for original trials for both metrics on Figure 19b). The clustering purity highlights a better performance of original trials over anonymized trials. The listeners obtain 86.40% clustering purity for original trials and 61.68% and 62.58% for both types of anonymized trials. These results indicate that linking an anonymized voice to its natural counterpart is not as easy as clustering unanonymized speakers. No significant difference between the two baselines is noticed for both metrics.

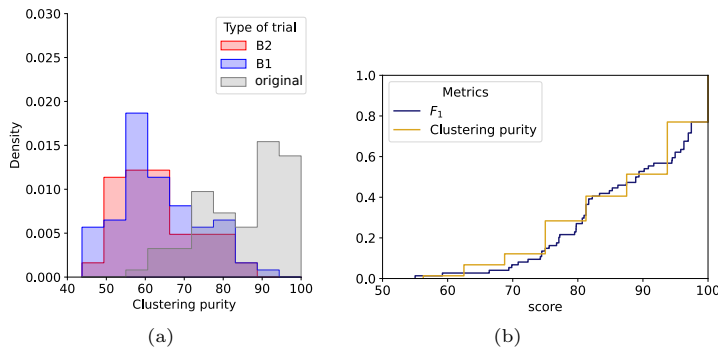


Figure 19: (a) density distributions for clustering purity; (b) cumulative density for clustering purity and  $F_1$  on the original (control) trials.

### 4.3. Comparison of objective and subjective evaluation results

In this section, we are interested in comparison of objective and subjective evaluation results. Given the subjective score on speaker verifiability (similarity), we computed EER,  $C_{\text{llr}}$ , and  $C_{\text{llr}}^{\text{min}}$ . We can then compare these metrics with those obtained based on objective evaluation. These comparisons are plotted in Figure 20 for EER.

The marker “Enr: o, Trl: o” in Figure 20 denotes unanonymized trials, and other markers denote anonymized trials from submitted systems. The comparison between unanonymized with anonymized trials indicates that both objective and subjective EERs increase after the trial is anonymized. However, the increase varies across the anonymization systems and test sets.

Furthermore, EERs of objective and subjective evaluation are positively correlated. It suggests that the concerned anonymization methods can hide the speaker identity to a moderate degree from both ASV system and human ears. This is an encouraging message from the challenge. Similar results can be observed for other privacy metrics:  $C_{\text{llr}}$ , and  $C_{\text{llr}}^{\text{min}}$  (see (Tomashenko et al., 2021), Section 5).

The relation between objective and subjective utility metrics is shown in Figure 21. The objective WERs have been reported in Figure 10. Subjective results on trials of female and male speakers were pooled, also pooled were the *common* and *different* subsets of *VCTK-test*. Results from the figure reinforced the observation that all involved anonymization systems degraded the objective and subjective utility metrics. For *LibriSpeech-test*, best results for all utility

metrics was achieved using a signal-based approach **I1**, and the worst — by **K2**. However, for *VCTK-test* there is no system that is the best (or worst) for all metrics. This is mostly due to the fact that objective utility evaluation results (decoding errors) are less consistent for different datasets than the subjective ones (naturalness and intelligibility).

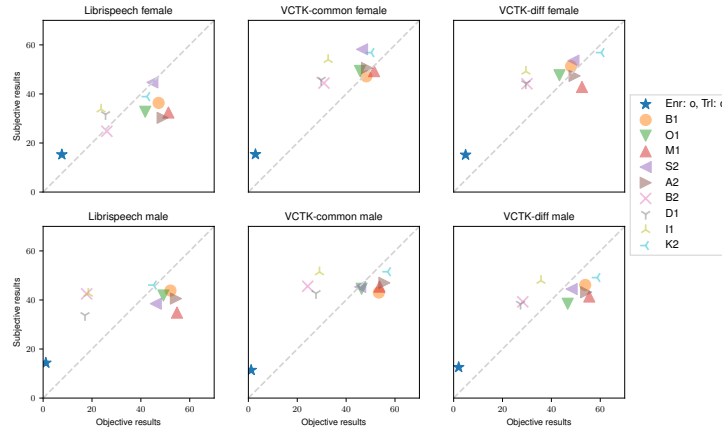


Figure 20: Objective versus subjective EER, %

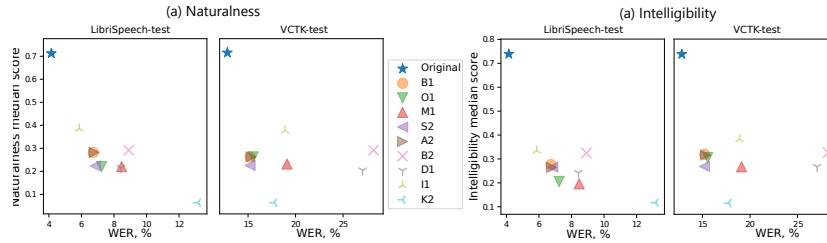


Figure 21: Objective versus subjective utility metrics

## 5. Conclusions

The VoicePrivacy 2020 Challenge was conceived to promote private-by-design and private-by-default speech technology and is the first evaluation campaign in voice anonymization. The voice anonymization task is defined as a game between users and attackers, with three possible attack models each corresponding to adversaries with different knowledge of the applied anonymization methods. The paper describes a full evaluation framework for the benchmarking of different anonymization solutions, including datasets, experimental protocols and metrics, as well as two open-source baseline anonymization solutions in addition to the comprehensive objective and subjective evaluation of both baseline



systems and those submitted by challenge participants. These indicate the potential for successful anonymization and serve as a platform for future work in what is now a burgeoning research field.

### 5.1. Summary and findings

The challenge attracted participants from both academia and industry, including from experts already working on anonymization as well as those new to the field. Anonymization systems designed by challenge participants are effective in terms of anonymization, each with different trade offs between privacy and utility. Submitted systems can be broadly classified in to two different classes of anonymization approaches: x-vectors-based methods with speech synthesis models (such as the primary baseline **B1**) and signal-processing based methods (relating to the secondary baseline **B2** and system **I1**). The x-vector based methods in average provide the best results in terms of objective evaluation.<sup>14</sup> In contrast, subjective evaluation shows that signal-processing based methods tend to give better results in terms of naturalness and intelligibility.

More consistent findings show that anonymization produced by all systems degrade naturalness and intelligibility, as well as the WER. Furthermore, the best systems in terms of WER are based on x-vector anonymization whereas the best system in terms of intelligibility is system **I1**.

Anonymization is also achieved only partially and always at the cost of utility; no single system gives the best performance for all metrics and each system offers a different trade-off between privacy and utility, whether judged objectively or subjectively. This finding holds no matter what the attack model. While for the ignorant attack model, many systems achieve EERs above 50%, for the lazy-ignorant attack model, best results are in the range of about 33 – 43% , and in the range of 16 – 26% for the semi-informed attack model. System rankings are also different in each case, demonstrating the challenge to design anonymisation systems that perform well across the range of different VoicePrivacy attacks models.

Challenge participants investigated the proposed anonymization approaches and suggested improvements in some test-cases over the baseline anonymization solutions. They found out, that (1) resynthesis alone degrades utility, while also improving privacy; (2) there is potential for privacy leakage not only in x-vector embeddings, but also in phonetic features and pitch estimates (Champion et al., 2020b; Mawalim et al., 2020); (3) the distribution of anonymized x-vectors differs to that of original x-vectors (Turner et al., 2020). Recent work shows the potential to reduce privacy leakage in pitch estimates while also protecting utility (Champion et al., 2020a; Srivastava et al., 2021). Other findings show that degradations to utility can be mitigated by retraining models used in downstream tasks, such as ASR, using anonymized data. Lastly, we identified

---

<sup>14</sup>There are some exceptions, in particular, related to the WER results for system **I1** and the *LibriSpeech* dataset

some differences or bias in performance across different datasets and for different speaker gender. The scale of these differences in one factor, among others discussed below, that warrants further attention in future research.

### 5.2. Open questions and future directions

A common understanding of VoicePrivacy is still in its infancy. For one, communicating the achieved in layperson terms remains a challenge to better integrate the larger speech community and for outreach to the public at large; for another, VoicePrivacy cannot remain at scratching the surface of privacy issues related to speech and language technology. In the first edition, while considering biometric identity as sensitive information, there are other types of sensitive information encoded and transported through speech as a communication medium. Moreover, by constraining the first edition to the operability of speech recognition, linguistic features still allow for extracting biometric characteristics to identify authorship. Depending on the context, the settings of ASV and ASR systems, one might argue that for prompted speech in automated call centers, there is less subjective variability in what is said; let alone, the goal of VoicePrivacy as a community is speech technology as a whole.

Future editions of the VoicePrivacy Challenge will include stronger baseline solutions, possible extensions of the tasks, and re-visited evaluation protocols:

- *Improved anonymization methods for stronger baseline solutions.* For the primary baseline and related approaches, perspective improvements in x-vector based anonymization include adversarial learning (Espinoza-Cuadros et al., 2020a) and design strategies based on speaker space analysis, gender, distance metric, etc. (Srivastava et al., 2020a, 2021). Sensitive information can be further removed from prosodic and other features, in particular, from pitch (Srivastava et al., 2021; Champion et al., 2020a) and phonetic (BN) features. Improved algorithms to use the speaker pool should take into account not only speaker characteristics before anonymization but also voice distinctiveness after anonymization. Moreover, the quality of the synthesized speech using unseen x-vectors has room for improvement. For the secondary baseline, we will consider its extension using a stochastic choice of McAdam’s coefficient (Patino et al., 2020).
- *Stronger and more realistic attack models.* Development and investigation of stronger attack models is another potential direction. A knowledgeable and experienced adversary will improve the ASV system and adapt it to make better decisions, i.e., to yield better class discrimination alongside accurate forecasts. Contrary to the conventional experimental validation based on error rates, an adversary actually needs to put a specific threshold and might want to change this threshold, depending on the settings of the ASV systems. In other words, priors and costs that determine the decision policy of an adversary need to be highly adaptable.
- *Alternative privacy and utility evaluation metrics.* The ongoing work on privacy preservation assessment is focusing on the development of new eval-

uation frameworks, anonymization metrics, and investigation of their correlation and complementarity. This includes the ZEBRA framework (Nautsch et al., 2020; Noé et al., 2021), objective and subjective linkability metrics (Maouche et al., 2020). Also one may be interested in evaluation that is close to real industry applications and tasks, for example, speaker labeling for diarization, analysis of time and quality required for annotation of real vs anonymized speech (Espinoza-Cuadros et al., 2020b). The metrics considered in the challenge do not evaluate fully the anonymization requirement that all characteristics in speech signal except the speaker identity, should be intact, and some speech characteristics (such as prosody) were not evaluated. Relevant utility metrics depend on the user’s downstream tasks, and for additional downstream tasks other utility metrics should be considered.

- *Attributes.* Besides the speaker identity information, speech also conveys other attributes that can be considered as sensitive, such as emotional state, age, gender, accent, etc. Selective suppression of such attributes is a possible task extension.
- *Privacy vs utility trade-off.* The privacy is often achieved at the expense of utility, and an important question is how to set up a proper threshold between privacy and utility (Li & Li, 2009). When developing anonymization methods, a joint optimization of utility gain and privacy loss can be performed by incorporating them into the criterion for training anonymization models (Kai et al., 2021).
- *Integrated approach to voice privacy and security.* In the bigger picture, security and privacy need to be thought of together and not as opposing forces: positive-sum solutions (Cavoukian, 2017) need to be sought to design technology for better products and services. In other words, while one might draw inspiration from machine learning, forensic sciences, and biometrics, integrated privacy designs for speech and language technology must sacrifice neither security, business interests, nor privacy. Developing of adequate VoicePrivacy safeguards demands future directions that empower capacity for their credible and adequate use in integrated privacy designs which beyond technology include organisational measures.

## Acknowledgment

VoicePrivacy was born at the crossroads of projects VoicePersonae, COMPRISE (<https://www.compriseh2020.eu/>), and DEEP-PRIVACY. Project HARPOCRATES was designed specifically to support it. The authors acknowledge support by ANR, JST (21K17775), and the European Union’s Horizon 2020 Research and Innovation Program, and they would like to thank Christine Meunier.

## References

- Aloufi, R., Haddadi, H., & Boyle, D. (2020). Privacy-preserving voice analysis via disentangled representations. In *Proceedings of the 2020 ACM SIGSAC Conference on Cloud Computing Security Workshop* (pp. 1–14).
- Bahmaninezhad, F., Zhang, C., & Hansen, J. H. (2018). Convolutional neural network based speaker de-identification. In *Odyssey* (pp. 255–260).
- Brümmer, N., & Du Preez, J. (2006). Application-independent evaluation of speaker detection. *Computer Speech and Language*, *20*, 230–275.
- Cavoukian, A. (2017). Global privacy and security, by design: Turning the “privacy vs. security” paradigm on its head. *Health and Technology, Privacy and Security of Medical Information*, *7*, 329–333.
- Champion, P., Jouvét, D., & Larcher, A. (2020a). *A Study of F0 Modification for X-Vector Based Speech Pseudo-Anonymization Across Gender*. Research Report INRIA Nancy, équipe Multispeech. URL: <https://hal.archives-ouvertes.fr/hal-02995862>.
- Champion, P., Jouvét, D., & Larcher, A. (2020b). *Speaker information modification in the VoicePrivacy 2020 toolchain*. Research Report INRIA Nancy, équipe Multispeech ; LIUM - Laboratoire d’Informatique de l’Université du Mans. URL: <https://hal.archives-ouvertes.fr/hal-02995855>.
- Chung, J. S., Nagrani, A., & Zisserman, A. (2018). VoxCeleb2: Deep speaker recognition. In *Interspeech* (pp. 1086–1090).
- Dubagunta, S. P., van Son, R. J., & Doss, M. M. (2020). Adjustable deterministic pseudonymisation of speech: Idiap-NKI’s submission to VoicePrivacy 2020 challenge, . URL: <https://www.voiceprivacychallenge.org/docs/Idiap-NKI.pdf>.
- Dwork, C. (2009). The differential privacy frontier. In *Theory of Cryptography Conference* (pp. 496–502). Springer.
- Espinoza-Cuadros, F. M., Perero-Codosero, J. M., Antón-Martín, J., & Hernández-Gómez, L. A. (2020a). Speaker de-identification system using autoencoders and adversarial training, . [arXiv:2011.04696](https://arxiv.org/abs/2011.04696).
- Espinoza-Cuadros, F. M., Perero-Codosero, J. M., Antón-Martín, J., & Hernández-Gómez, L. A. (2020b). Speaker de-identification system using autoencoders and adversarial training. *Presentation at the VoicePrivacy 2020 virtual workshop*: <https://www.voiceprivacychallenge.org>, . URL: <https://youtu.be/wCvIh4G3fFM>.
- Fang, F., Wang, X., Yamagishi, J., Echizen, I., Todisco, M., Evans, N., & Bonastre, J.-F. (2019). Speaker anonymization using x-vector and neural waveform models. In *Speech Synthesis Workshop* (pp. 155–160). doi:10.21437/SSW.2019-28.

- Ghorshi, S., Vaseghi, S., & Yan, Q. (2008). Cross-entropic comparison of formants of British, Australian and American English accents. *Speech Communication*, 50, 564–579.
- Griffin, D., & Lim, J. (1984). Signal estimation from modified short-time fourier transform. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 32, 236–243.
- Gupta, P., Prajapati, G. P., Singh, S., Kamble, M. R., & Patil, H. A. (2020). Design of voice privacy system using linear prediction. URL: <https://www.voiceprivacychallenge.org/docs/DA-IICT-Speech-Group.pdf>.
- Han, Y., Li, S., Cao, Y., Ma, Q., & Yoshikawa, M. (2020a). Voice-indistinguishability: Protecting voiceprint in privacy-preserving speech data release. *arXiv preprint arXiv:2004.07442*, .
- Han, Y., Li, S., Cao, Y., & Yoshikawa, M. (2020b). System description for voice privacy challenge. kyoto team, . URL: <https://www.voiceprivacychallenge.org/docs/Kyoto.pdf>.
- Hashimoto, K., Yamagishi, J., & Echizen, I. (2016). Privacy-preserving sound to degrade automatic speaker verification performance. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 5500–5504).
- Hintze, J. L., & Nelson, R. D. (1998). Violin plots: a box plot-density trace synergism. *The American Statistician*, 52, 181–184.
- Huang, C.-L. (2020). Analysis of pingan submission in the voiceprivacy 2020 challenge. URL: <https://www.voiceprivacychallenge.org/docs/PingAn.pdf>.
- Jin, Q., Toth, A. R., Schultz, T., & Black, A. W. (2009). Speaker de-identification via voice transformation. In *ASRU*.
- Kai, H., Takamichi, S., Shiota, S., & Kiya, H. (2021). Lightweight voice anonymization based on data-driven optimization of cascaded voice modification modules. In *IEEE SLT 2021* (pp. 560–566). IEEE.
- Li, T., & Li, N. (2009). On the tradeoff between privacy and utility in data publishing. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining* (pp. 517–526).
- Magariños, C., Lopez-Otero, P., Docio-Fernandez, L., Rodriguez-Banga, E., Erro, D. et al. (2017). Reversible speaker de-identification using pre-trained transformation functions. *Computer Speech & Language*, 46, 36–52.
- Maouche, M., Srivastava, B. M. L., Vauquier, N., Bellet, A., Tommasi, M., & Vincent, E. (2020). A Comparative Study of Speech Anonymization Metrics. In *Proc. Interspeech 2020* (pp. 1708–1712). URL: <http://dx.doi.org/10.21437/Interspeech.2020-2248>. doi:10.21437/Interspeech.2020-2248.

- Mawalim, C. O., Galajit, K., Karnjana, J., & Unoki, M. (2020). X-Vector Singular Value Modification and Statistical-Based Decomposition with Ensemble Regression Modeling for Speaker Anonymization System. In *Proc. Interspeech 2020* (pp. 1703–1707). URL: <http://dx.doi.org/10.21437/Interspeech.2020-1887>. doi:10.21437/Interspeech.2020-1887.
- McAdams, S. (1984). Spectral fusion, spectral parsing and the formation of the auditory image. *Ph. D. Thesis, Stanford, .*
- Morise, M., Yokomori, F., & Ozawa, K. (2016). World: a vocoder-based high-quality speech synthesis system for real-time applications. *IEICE TRANSACTIONS on Information and Systems, 99*, 1877–1884.
- Moulines, E., & Charpentier, F. (1990). Pitch-synchronous waveform processing techniques for text-to-speech synthesis using diphones. *Speech communication, 9*, 453–467.
- Nagrani, A., Chung, J. S., & Zisserman, A. (2017). VoxCeleb: a large-scale speaker identification dataset. In *Interspeech* (pp. 2616–2620).
- Nautsch, A., Jasserand, C., Kindt, E., Todisco, M., Trancoso, I., & Evans, N. (2019a). The GDPR & speech data: Reflections of legal and technology communities, first steps towards a common understanding. In *Interspeech* (pp. 3695–3699).
- Nautsch, A., Jimenez, A., Treiber, A., Kolberg, J., Jasserand, C., Kindt, E., Delgado, H. et al. (2019b). Preserving privacy in speaker and speech characterisation. *Computer Speech and Language, 58*, 441–480.
- Nautsch, A., Patino, J., Tomashenko, N., Yamagishi, J., Noé, P.-G., Bonastre, J.-F., Todisco, M., & Evans, N. (2020). The Privacy ZEBRA: Zero Evidence Biometric Recognition Assessment. In *Proc. Interspeech 2020* (pp. 1698–1702). URL: <http://dx.doi.org/10.21437/Interspeech.2020-1815>. doi:10.21437/Interspeech.2020-1815.
- Noé, P.-G., Bonastre, J.-F., Matrouf, D., Tomashenko, N., Nautsch, A., & Evans, N. (2020). Speech Pseudonymisation Assessment Using Voice Similarity Matrices. In *Proc. Interspeech 2020* (pp. 1718–1722). URL: <http://dx.doi.org/10.21437/Interspeech.2020-2720>. doi:10.21437/Interspeech.2020-2720.
- Noé, P.-G., Nautsch, A., Evans, N., Patino, J., Bonastre, J.-F., Tomashenko, N., & Matrouf, D. (2021). Two assessment frameworks for privacy preserving transformation of speech. Submitted.
- O’Brien, B., Tomashenko, N., Chanclu, A., & Bonastre, J.-F. (2021). Anonymous speaker clusters: Making distinctions between anonymised speech recordings with clustering interface. In *Interspeech 2021*.

- Panayotov, V., Chen, G., Povey, D., & Khudanpur, S. (2015). LibriSpeech: an ASR corpus based on public domain audio books. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 5206–5210).
- Patino, J., Todisco, M., Nautsch, A., & Evans, N. (2020). *Speaker anonymisation using the McAdams coefficient*. Technical Report EURECOM+6190 Eurecom. URL: <http://www.eurecom.fr/publication/6190>.
- Patino, J., Tomashenko, N., Todisco, M., Nautsch, A., & Evans, N. (2020). Speaker anonymisation using the McAdams coefficient. *arXiv:2011.01130*.
- Peddinti, V., Povey, D., & Khudanpur, S. (2015). A time delay neural network architecture for efficient modeling of long temporal contexts. In *Interspeech* (pp. 3214–3218).
- Pobar, M., & Ipšić, I. (2014). Online speaker de-identification using voice transformation. In *37th International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO)* (pp. 1264–1267).
- Povey, D., Cheng, G., Wang, Y., Li, K., Xu, H., Yarmohammadi, M. et al. (2018). Semi-orthogonal low-rank matrix factorization for deep neural networks. In *Interspeech* (pp. 3743–3747).
- Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N. et al. (2011). The Kaldi speech recognition toolkit.
- Qian, J., Du, H., Hou, J., Chen, L., Jung, T., Li, X.-Y., Wang, Y., & Deng, Y. (2017). Voicemask: Anonymize and sanitize voice input on mobile devices. *arXiv preprint arXiv:1711.11460*, .
- Qian, J., Han, F., Hou, J., Zhang, C., Wang, Y., & Li, X.-Y. (2018). Towards privacy-preserving speech data publishing. In *2018 IEEE Conference on Computer Communications (INFOCOM)* (pp. 1079–1087).
- Ramos, D., & Gonzalez-Rodriguez, J. (2008). Cross-entropy analysis of the information in forensic speaker recognition. In *Odyssey*.
- Rosenberg, A., & Ramabhadran, B. (2017). Bias and Statistical Significance in Evaluating Speech Synthesis with Mean Opinion Scores. In *Proc. Interspeech* (pp. 3976–3980).
- Snyder, D., Garcia-Romero, D., Sell, G., Povey, D., & Khudanpur, S. (2018). X-vectors: Robust DNN embeddings for speaker recognition. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 5329–5333).
- van Son, R. (2020). Adjustable Deterministic Pseudonymization of Speech Listening Experiment, Report of listening experiments. URL: <https://doi.org/10.5281/zenodo.3773931>. doi:10.5281/zenodo.3773931.

- Srivastava, B. M. L., Bellet, A., Tommasi, M., & Vincent, E. (2019). Privacy-preserving adversarial representation learning in ASR: Reality or illusion? In *Interspeech* (pp. 3700–3704).
- Srivastava, B. M. L., Maouche, M., Sahidullah, M., Vincent, E., Bellet, A., Tommasi, M., Tomashenko, N., Wang, X., & Yamagishi, J. (2021). On the robustness of speaker anonymization from speaker, user and attacker’s perspective. *submitted*, .
- Srivastava, B. M. L., Tomashenko, N., Wang, X., Vincent, E., Yamagishi, J., Maouche, M., Bellet, A., & Tommasi, M. (2020a). Design choices for x-vector based speaker anonymization. In *Interspeech*.
- Srivastava, B. M. L., Vauquier, N., Sahidullah, M., Bellet, A., Tommasi, M., & Vincent, E. (2020b). Evaluating voice conversion-based privacy protection against informed attackers. In *2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 2802–2806).
- Tomashenko, N., Srivastava, B. M. L., Wang, X., Vincent, E., Nautsch, A., Yamagishi, J., Evans, N., Patino, J., Bonastre, J.-F., Noé, P.-G., & Todisco, M. (2020a). Introducing the VoicePrivacy Initiative. In *Proc. Interspeech 2020* (pp. 1693–1697). URL: <http://dx.doi.org/10.21437/Interspeech.2020-1333>. doi:10.21437/Interspeech.2020-1333.
- Tomashenko, N., Srivastava, B. M. L., Wang, X., Vincent, E., Nautsch, A., Yamagishi, J., Evans, N. et al. (2020b). Post-evaluation analysis for the VoicePrivacy 2020 challenge: Using anonymized speech data to train attack models and ASR, . URL: [https://www.voiceprivacychallenge.org/docs/VoicePrivacy2020\\_post\\_evaluation.pdf](https://www.voiceprivacychallenge.org/docs/VoicePrivacy2020_post_evaluation.pdf).
- Tomashenko, N., Srivastava, B. M. L., Wang, X., Vincent, E., Nautsch, A., Yamagishi, J., Evans, N. et al. (2020c). The VoicePrivacy 2020 Challenge evaluation plan, . URL: [https://www.voiceprivacychallenge.org/docs/VoicePrivacy\\_2020\\_Eval\\_Plan\\_v1\\_3.pdf](https://www.voiceprivacychallenge.org/docs/VoicePrivacy_2020_Eval_Plan_v1_3.pdf).
- Tomashenko, N., Wang, X., Vincent, E., Patino, J., Srivastava, B. M. L., Noé, P.-G., Nautsch, A., Evans, N., Yamagishi, J., O’Brien, B., Chanclu, A., Bonastre, J.-F., Todisco, M., & Maouche, M. (2021). Supplementary material to the paper. The VoicePrivacy 2020 Challenge: Results and findings. Submitted to Special Issue on Voice Privacy in Computer Speech and Language.
- Turner, H., Lovisotto, G., & Martinovic, I. (2020). Speaker anonymization with distribution-preserving x-vector generation for the VoicePrivacy Challenge 2020, . [arXiv:2010.13457](https://arxiv.org/abs/2010.13457).
- Veaux, C., Yamagishi, J., & MacDonald, K. (2019). CSTR VCTK corpus: English multi-speaker corpus for CSTR voice cloning toolkit (version 0.92). URL: <https://datashare.is.ed.ac.uk/handle/10283/3443>.



- Vincent, E., Nautsch, A., Evans, N., Tomashenko, N., Yamagishi, J., Smaragdis, P., & Bonastre, J.-F. (2021). An Overview of Speech Privacy Research. Submitted to Special Issue on Voice Privacy in Computer Speech and Language.
- Wang, X., & Yamagishi, J. (2019). Neural harmonic-plus-noise waveform model with trainable maximum voice frequency for text-to-speech synthesis. In *Speech Synthesis Workshop* (pp. 1–6). doi:10.21437/SSW.2019-1.
- Watanabe, S., Hori, T., Kim, S., Hershey, J. R., & Hayashi, T. (2017). Hybrid ctc/attention architecture for end-to-end speech recognition. *IEEE Journal of Selected Topics in Signal Processing*, 11, 1240–1253.
- Yoo, I.-C., Lee, K., Leem, S., Oh, H., Ko, B., & Yook, D. (2020). Speaker anonymization for personal information protection using voice conversion techniques. *IEEE Access*, 8, 198637–198645.
- Zen, H., Dang, V., Clark, R., Zhang, Y., Weiss, R. J., Jia, Y., Chen, Z., & Wu, Y. (2019). LibriTTS: A corpus derived from LibriSpeech for text-to-speech. In *Interspeech* (pp. 1526–1530).