

# Adjustable Deterministic Pseudonymization of Speech

Rob J.J.H. van Son<sup>1,2</sup>

<sup>1</sup>NKI-AVL, Amsterdam; <sup>2</sup>ACLIC, University of Amsterdam, The Netherlands

r.v.son@nki.nl

## Abstract

A method for the pseudonymization of speech is presented that allows to obfuscate the identity of a speaker in untranscribed running speech. The approach is to manipulate the spectro-temporal structure of the speech to simulate a different length and structure of the vocal tract, as well as a different pitch and speaking rate. The method is deterministic, and partially reversible. The extend of the changes is adjustable and gradual. A series of ABX listening experiments show that both experts and non-experts identify speakers in less than 70% of forced choice pairs while listeners are able to identify over 90% of speakers without pseudonymization. Reverting the procedure, de-pseudonymization, is partially effective. Some pseudonymization targets, e.g., those simulating a long vowel tract, are more amenable to de-pseudonymization than others. The method also works differently on female and male voices. Depending on the pseudonymization target, female speakers were less well identified after pseudonymization and de-pseudonymization than male speakers.

**Index Terms:** speech pseudonymization, human speaker identification,

## 1. Introduction

Collecting and sharing speech recordings is important for progress in speech science and technology. A lot of progress in speech technology has been made possible by the availability of large speech corpora in combination with advanced statistical techniques [1, 2]. Speech recordings are also a possible privacy risks. This is especially true when the speakers have medical conditions, are minors, or the subject matter is sensitive. But these are also groups that might benefit from improvements in speech technology tailored to their needs. The privacy risks resulting from sharing speech recordings would be mitigated if the probability of speaker (re-)identification could be reduced while retaining useful linguistic and para-linguistic features. Such a pseudonymization of speech would shift the risk-benefit balance for sharing speech corpora towards more sharing. However, current state-of-the art speech pseudonymization methods are not good at preserving the linguistic and para-linguistic features of interest, either because the transformed speech is degraded [3] or because the connection with the source speech is cut by re-synthesizing the speech in another voice using chained ASR and speech synthesis [4]. The research community has acknowledged these problems and a call for improving pseudonymization of speech has gone out [5].

The literature on data anonymization (e.g., [6, 7]) can be crudely summarized as “anonymous data is not useful, useful data is not anonymous”. This is also likely to be true for speech transforms. Contrary to anonymization, which is either *on* or *off*, pseudonymization is based on the assumption that data can be re-identified, in principle, but additional information is needed to do so. The risk of re-identification of pseudonymized

data is then the risk that the missing information can be reconstructed by an attacker.

Pseudonymization of speech will always be a trade-off between risk of re-identification and usefulness. This suggests an approach to pseudonymization that is adjustable in the level of information removed from the speech while still preserving relevant features enough to make the result useful. The aim is to develop a method in which the transformation of the speech can be tailored to the risk profile and features needed (c.f., [8]).

Two sources of speaker variation that are useful for speaker identification can be distinguished, *inherent* features, i.e., those that derive from a speaker’s anatomy and physiology, and *learned* features [9]. This study aims at hiding global, inherent features of speakers, more specific, vocal tract related spectral features (c.f., [10]). This translates to a method to make changes in speech that relate to vocal tract length, average formant frequencies and intensities, pitch, and speaking rate. The information that is hidden are the original values of these measures and the extend of the changes. The current study tries to produce high quality, pseudonymized speech. The risk of re-identification is investigated in ABX listening experiments using human listeners. In cases where privacy risks should be reduced further, the target speech can be degraded more, e.g., by removing or randomizing spectral bands. This approach is not investigated here.

The pseudonymization procedure is discussed in section 2.1. The listening experiments described in section 2.2. The results are presented in section 3 and discussed in section 4.

## 2. Methods

### 2.1. Pseudonymization

The method for the pseudonymization of speech used in the current study is based on two basic procedures (software available on GitHub [11]). To change the perceived acoustic length of the vocal tract, the playback speed of the sound is changed. Technically, this is achieved by changing the sampling rate for playback. Overlap-add [12] is used to adjust the pitch and the duration of the utterances. Vocal tract length corresponds to formant values in that an increase of vocal tract length by a factor  $\alpha$  induces a formant shift by a factor  $1/\alpha$ . In the remainder of this paper, the estimated vocal tract length will be represented by the neutral first resonance frequency  $\phi$ .

The program *Praat* has two commands that both perform these two operations using the same algorithms: *Change gender...* and *Change speaker...* This study uses the *Change gender...* command internally because the command options better suited the current approach. In these commands, the desired new pitch median is set as an absolute value, but it depends on correct pitch measurement in the source speech. Both commands work on the vocal tract length and duration by a *Formant shift ratio* and a *Duration factor*. To implement a change to a specified target vocal tract length and duration, or speaking rate,

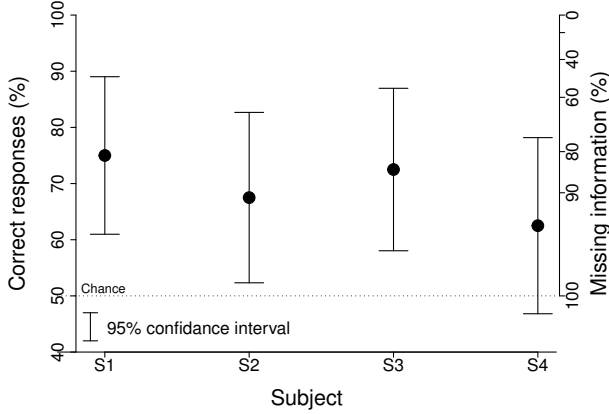


Figure 1: *Speaker identification in experiment 1 by expert subject with correct responses (left) and missing information (right, 100% = 1 bit). Confidence intervals from Student distribution. Overall mean correct: 69%, 95% conf int. [61, 78]%. No differences were found in responses to male and female speakers.*

the estimated vocal tract length and speaking rate of the source speaker have to be supplied.

The speaker’s pitch, vocal tract length and structure, and “typical” speaking rate are estimated from a collection of untranscribed speech recordings, preferably over 300 seconds of speech spoken in a comparable style. A speaker’s  $\hat{\phi}$  is estimated from the first 4 formant frequencies according to eq. 20 of [13] using the proposed extension (table 3, *ibid.*):

$$\hat{\phi} = 229 + 0.030F_1/1 + 0.082F_2/3 + 0.124F_3/5 + 0.354F_4/7 \quad (1)$$

Vocal tract length (VTL) is the mean VTL calculated from formant values of points closest to ( $F_1=500$ ,  $F_2=1500$ ), determined using the *Praat* `robust` formant option [14]. Vowels are segmented using the method of [15]. Speaking rate is determined by the syllable rate determined from a modified version of a script by De Jong and Wempe [16] taken from [15].

Speakers do not only differ in vocal tract length, but also in the vocal tract structure that changes the global (median) position and width and height of formants. Individual formant positions and intensities are changed by first creating a version of the speech utterance with the median formant at the right frequency and intensity using the above procedure, using  $\hat{\phi}_i = \text{median}(F_i/(2i-1))$  where  $i$  is the number of the target formant. The final sound is constructed by replacing the target bands in the main sound by the corresponding band in the adapted sound (using stop and pass Hann band filters  $F_i \pm \phi$  with smoothing  $(F_i - \phi)/10$ ). In addition to bands around the formants, the lower part of the spectrum is treated as a separate band ( $F_0 : 0 - \phi/2$ ). The speaker frequencies and intensities are estimated as the median formant values and intensities in the pass bands of a sound normalized to 70 dB.

Frequencies can be randomly chosen in the range  $F_i \pm 40$  and  $F_i \pm 75$  Hz for  $F_{0-1}$  and  $F_{2-5}$ , respectively. Intensities can be randomly chosen in the range  $64 \pm 4.5$ ,  $67 \pm 2.5$ ,  $58 \pm 4.5$ ,  $50 \pm 8$ ,  $47 \pm 10$ ,  $45 \pm 9$  dB ( $F_{0-5} \pm 2SD$ ). For the pseudonymization in the current study, the positions (frequencies) and intensities of the bands  $F_0$ , and  $F_{3-5}$  were randomly set. Alternatively, the profile values of a selected speaker can be used as targets (experiment 3, see below).

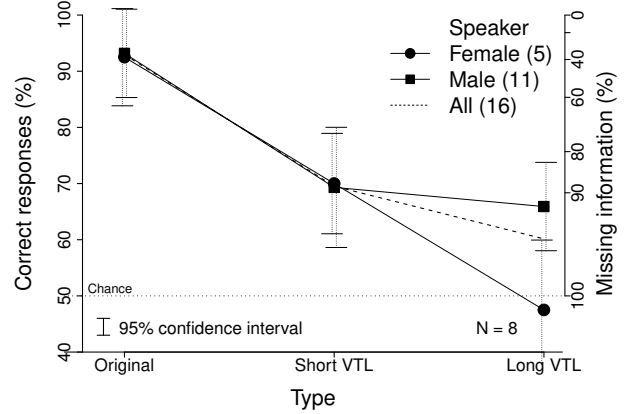


Figure 2: *Speaker identification in experiment 2 by stimulus type and speaker gender. Original: AB are original recordings, Short VTL: AB pseudonymized to a short vocal tract length, Long VTL: AB pseudonymized to a long vocal tract length. N: Number of subjects. See also Fig. 1.*

## 2.2. Listening experiments

Pseudonymized sentences and sentence fragments were produced by running the `PseudonymizeSpeech.praat` script with target values for the Long Vocal Tract (LVT):  $\phi = 500\text{Hz}$ ,  $F_0 = 120\text{Hz}$ ; and the Short Vocal Tract (SVT):  $\phi = 575\text{Hz}$ ,  $F_0 = 175\text{Hz}$ . Randomized values were used for the frequencies and intensities of bands  $F_0$ ,  $F_{3-5}$  (see above). Three experiments were performed. In Experiment 1, the target speaking rate was 3.8 and 4.2 syll/s (SVT and LVT), in experiments 2 and 3 the rate was 4.0 syll/s. Experiments are available at [17].

Listeners were presented with stimuli in three self-paced, online ABX forced choice experiments [18], where either stimulus *A* or *B* was from the same speaker as the target stimulus *X* (see Table 1). Subjects could listen to the three stimuli in any order and as often as they wanted. The position of the target speaker in *A* or *B* was randomized in the list. This position was identical for all listeners. Items were presented in pseudo-random order, different for each listener. Each experiment started with 4 practice items which were the last 4 items in reverse order.

**Experiment 1: Pseudosentences** from the IFA corpus read by 10 Dutch speakers (5F) [19]. Speaker profiles were derived from all pseudosentences read by that speaker. LVT and SVT examples of the target speaker and detractor, were presented to 4 experts, 3 speech therapists and 1 linguist. When *X* was LVT, *A* and *B* were SVT and vice versa. Each target speaker was presented once with a male and once with a female detractor.

**Experiment 2: Sentence fragments** with a maximum duration of 3s were selected from readings of *Treasure Island* taken from the *Parallel Audiobook Corpus* [20] read by 16 speakers of British English (5F). Speaker profiles were derived from all sentences in a single chapter, not used for selecting stimulus sentences. *X* was an original recording from the target speaker, *A* and *B* were both either Original recordings, or LVT or SVT pseudonymizations. There were 16 ABX stimulus combinations for each condition, Original, LVT, and SVT, 48 ABX combinations in total. Each speaker was used only once as target speaker for each condition (not counting practice items). Detractors were selected at random irrespective of gender. The genders of target speaker and detractor were the same (ff or mm)

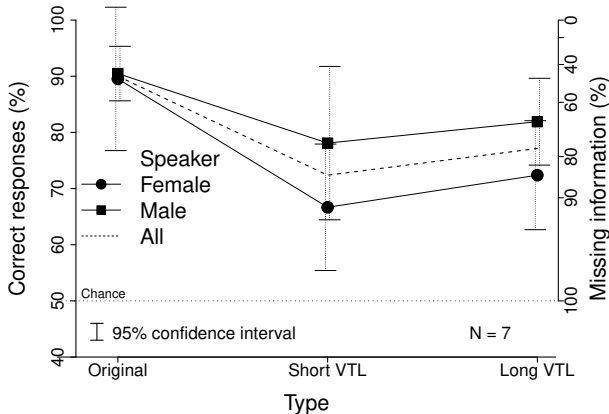


Figure 3: Identification after de-pseudonymization in experiment 3 by stimulus type and speaker gender. Speaker: Target speakers, 15F/15M for each Type, 90 in total. See also Fig. 2.

for 27 stimuli and different (fm or mf) for 21 stimuli. For this experiment, 8 “naive” listeners participated, recruited by email.

**Experiment 3:** All sentences from the *Bonafide* recordings from the Logical Access part of the 2019 ASV spoof corpus [21] were pseudonymized as in experiment 2. Speaker profiles were derived from all sentences of that speaker. Gender information was available for 58 out of 107 speakers. A linear model based on the speaker profiles, with perfect fit on the known genders, was used to predict the gender of the other speakers. Sentence fragments with a maximum duration of 3s were selected as ABX stimuli from target speaker and detractor, and were all back-transformed using the speaker profile of the target speaker. In the back-transform, the formant frequencies and band intensities of the transformed segments were not known. Therefore, only the vocal tract length, pitch, and speaking rate were transformed to the target speaker profile. Target speaker and detractor were always of the same gender, both male or female. This was done because pilot tests showed that mixed gender stimulus pairs were perfectly identified. Each condition in experiment 3, Original, LVT, and SVT, contained sentences from 15 male and 15 female target speakers and randomly selected detractors of the same gender, 90 ABX stimulus sets in total. In experiment 3, each speaker was only used once as target speaker and once as detractor (not counting practice items). Sentences were selected at random from the corpus from each speaker, but no sentence recording was used twice in the experiment.

Subjects for experiments 2 and 3 were recruited over email with written instructions. Listening conditions in these two experiments were not supervised. As quality assurance, only responses from subjects who were able to correctly identify 70% of the target speakers in the original recordings (condition *Original*) were included in the analysis. Five subjects participated in both experiment 2 and 3, one in experiments 1 and 3.

Table 1: Summary of ABX listening experiments. Sp.: Speakers.

Exp	Corpus	Speech ( $\leq 3s$ )	Sp. F/M	Subjects
1	[19]	Pseudo sent.	5/5	4 experts
2	[20]	sentences	5/11	8 naive
3	[21]	sentences	45/45	7 naive

No formal evaluation of stimulus quality was performed. However, informal impressions of stimulus quality in experiments 1 and 2 were from audibly distorted but highly intelligible to near natural. In experiment 3, the second de-pseudonymizing transform did markedly affect the speech quality. Short and Long VTL stimuli in experiment 3 were still intelligible, but sounded clearly distorted.

### 3. Results

All statistical analysis is done with R [22]. Missing information is calculated as the entropy  $H = -\sum_{i=1}^2 p_i \log_2 p_i$  (in %). Differences in identification between conditions and stimulus classes are tested using paired Student t-tests (following [23]).

#### 3.1. Listening experiments

**Experiment 1:** The expert listeners reported that this was a difficult task where they found it difficult to believe that the target speaker was always among the response choices. The expert listeners identified the target speaker approximately 70% of the time (see Fig. 1, missing  $>80\%$  of information  $H$ ). The responses were better than chance and worse than perfect ( $p \leq 0.006$  for both 90% and 50% correct, t-test, not shown). There were no statistically significant differences between listeners and no effects of speaker gender (not shown).

**Experiment 2:** Responses of one subject, who did not reach 70% correct identification on the original recordings, were dropped (subject removed). On average, speaker identification of the original recordings was over 90% correct (see Fig. 2). The naive listeners identified the target speaker approximately 70% of the time in the short VTL condition and somewhat less in the long VTL condition (missing  $>80\%$  of information to identify the speaker). This was significantly less than in the original condition with unaltered speech ( $p \leq 0.0001$ , paired t-test by subject). The difference between the short and long VTL condition were not significant ( $p > 0.05$ ). There is a tentative difference in responses to the (5) female speakers and the (11) male speakers for the Long VT stimuli ( $p = 0.027$ , paired t-test). It appears that female speakers are not identified above chance level in the long VTL condition.

In the responses from experiments 1 and 2, there is a tendency that comparison to a detractor of a different gender improves identification of the target speaker (not shown). However, partly due to the design of the experiments, this could not be verified ( $p > 0.05$ , paired t-test).

**Experiment 3:** All subjects cleared the 70% correct criterion for the *Original* stimulus condition. Speaker identification of the original recordings in experiment 3 was around 90% correct (see Fig.3 and Table 2). De-pseudonymization, the in-

Table 2: Speaker identification accuracy in experiments 2 and 3. Linear mixed effects models of influence of (de-) pseudonymization and speaker gender on identification for each stimulus type (see text). Ex: Experiment,  $p$  (Ex):  $p$  value of difference between experiments,  $p$  (Gen):  $p$  value of difference between speaker genders in combined experiment.

Stimulus	Ex. 2 (sd)	Ex. 3 (sd)	$p$ (Ex)	$p$ (Gen)
Original	93% ( 6)	90% ( 8)	$>0.05$	$>0.05$
Short VTL	70% (11)	73% (12)	$>0.05$	$>0.05$
Long VTL	60% ( 7)	77% ( 6)	0.009	0.012

verse transform, was effective in reversing the pseudonymization towards a *Long VTL* target, increasing the identification from 60% to 78% correct (Table 2) with missing information  $\leq 80\%$  (Fig.3). However, the differences in identification between the *Original* and the de-pseudonymized stimuli was still significant ( $p \leq 0.009$  paired t-test by subject). The difference in identification between male and female was not significant in experiment 3 ( $p > 0.05$  for all stimulus types).

### 3.2. Modeling responses to experiments 2 & 3

The results of experiments 2 and 3 were combined in a linear mixed effect model to estimate the effects of speaker Gender and pseudonymization versus de-pseudonymization (Exp) on speaker Identification (I) for each stimulus type, i.e., *Original*, *Short VTL*, *Long VTL*. The full model was:

$$I \sim \text{Exp} + \text{Gender} + (\text{Exp} + \text{Gender} | \text{Subject}) \quad (2)$$

Subjects that participated in both experiments were identified in the model. Statistical significance was determined using ANOVA on full model versus a model with the relevant fixed factor removed. No difference was found for the *Original* and *Short VTL* stimuli ( $p > 0.05$ ). For the *Long VTL* target pseudonymizations, both the differences between male and female speakers and the differences between the experiments were statistically significant (see Table 2). Using the model of Eq. 2, male speakers were identified 13% points more than female speakers and de-pseudonymization increased identification by 19% points ( $p$  values in Table 2).

Experiment 3 only contained same gender comparisons between target and detractor speakers, while experiment 2 contained same and mixed gender comparisons. Same gender comparisons could be seen as “more difficult” than mixed gender comparisons. Repeating the linear mixed effect modelling with only the responses to same gender detractors gave the same results, no effect for *Original* and *Short VTL* stimuli ( $p > 0.05$ ) and a consistent effect for de-pseudonymization and speaker gender for *Long VTL* stimuli ( $p(\text{Ex}) = 0.008$ ,  $p(\text{Gen}) = 0.024$ , not shown). But the effect of de-pseudonymization only increases marginally (to 22% points). The overall effect of de-pseudonymization was found for both female and male speakers separately ( $p \leq 0.012$ , ANOVA, removing *Gender* from Eq. 2, not shown).

## 4. Discussion

All three ABX listening experiments showed reduced speaker identification after pseudonymization (Fig. 1 and 2) and also after de-pseudonymization (Fig. 3). After pseudonymization, more than 80% of the information necessary to make the choice between speaker A and B is lost ( $< 70\%$  correct identification, Fig. 1 and 2), compared to less than 40% information loss with the original recordings ( $> 90\%$  identification, Table 2). Reverting the pseudonymization transformation from known target positions can improve recognition, especially for speech transformed to a *Long VTL* (Fig. 3 and Table 2).

The responses in both experiments 2 and 3 displayed an asymmetry between male and female voices. Female speakers were identified worse than male speakers after pseudonymization as well as de-pseudonymization. This difference was statistically significant for the *Long VTL* condition when the responses in these experiments are combined (Table 2). This asymmetry was smaller, or absent, in the *Short VTL* condition (statistically not significant).

One objective of the current approach is to make pseudonymization deterministic and adjustable, i.e., gradual, on untranscribed recordings. The application described has all these features. It works in the spectro-temporal domain on any speech recording, and is intrinsically deterministic and reversible. However, reversibility is not necessarily an advantage. The exception to reversibility is the overlap-add procedure to adapt the pitch and duration of the speech which is inherently “lossy”, i.e., partially irreversible. But overlap-add is a well known, predictable, speech synthesis procedure. The aspects of the speech that are transformed as well as the extend of the changes can all be freely chosen. The only constraint is the quality of the resulting speech.

The primary aim of the pseudonymization is protection of the privacy of the speakers. Whether or not the levels of protection are sufficient depends on the needs of the situation and the risks that identification would pose. It is clear that the ability to, partially, de-pseudonymize the speech warrants extra attention. The current study explores one specific de-pseudonymization approach based on knowing the original pseudonymization target. An obvious way to hamper this de-pseudonymization approach would be to obfuscate this target by randomly selecting a pseudonymization target for each speaker, as estimating the target from short stretches of speech is difficult.

An important goal of pseudonymization of speech is to allow linguistic studies on speech samples without jeopardizing the privacy of the speakers. It is not yet known what linguistic and para-linguistic aspects of speech can still be studied after using the pseudonymization method described here. The current study only investigates the usefulness of pseudonymization for human listeners. How this method performs in combination with automatic speaker verification/identification applications will be investigated in another study.

## 5. Conclusions

A method to pseudonymize speech is described that is both deterministic and adjustable. The method can pseudonymize speech samples with only a few hundred seconds of speech of the source speaker. Pseudonymized samples are largely unidentifiable for human listeners. However, the deterministic nature of the procedure compel caution in applying the procedure to counter de-pseudonymization.

## 6. Acknowledgements

The Department of Head and Neck Oncology and Surgery of the Netherlands Cancer Institute receives a research grant from Atos Medical (Malmö, Sweden), which contributes to the existing infrastructure for quality of life research.

## 7. References

- [1] Z. Zhang, N. Cummins, and B. Schuller, "Advanced data exploitation in speech analysis: An overview," *IEEE Signal Processing Magazine*, vol. 34, no. 4, pp. 107–129, 2017.
- [2] Y. Ning, S. He, Z. Wu, C. Xing, and L.-J. Zhang, "A Review of Deep Learning Based Speech Synthesis," *Applied Sciences*, vol. 9, no. 19, p. 4050, Sep. 2019. [Online]. Available: <https://www.mdpi.com/2076-3417/9/19/4050>
- [3] B. M. L. Srivastava, N. Vauquier, M. Sahidullah, A. Bellet, M. Tommasi, and E. Vincent, "Evaluating Voice Conversion-based Privacy Protection against Informed Attackers," 2019. [Online]. Available: <https://hal.inria.fr/hal-02355115>
- [4] F. Fang, X. Wang, J. Yamagishi, I. Echizen, M. Todisco, N. Evans, and J.-F. Bonastre, "Speaker Anonymization Using X-vector and Neural Waveform Models," in *10th ISCA Speech Synthesis Workshop*. ISCA, Sep. 2019, pp. 155–160. [Online]. Available: [http://www.isca-speech.org/archive/SSW\\_2019/abstracts/SSW10.P.2-4.html](http://www.isca-speech.org/archive/SSW_2019/abstracts/SSW10.P.2-4.html)
- [5] N. Tomashenko, B. M. L. Srivastava, X. Wang, E. Vincent, A. Nautsch, J. Yamagishi, N. Evans, J.-F. Bonastre, P.-G. No  , M. Todisco, and J. Patino. (2020, Feb.) The VoicePrivacy 2020 Challenge. VoicePrivacy. [Online]. Available: <https://www.voiceprivacychallenge.org/>
- [6] I. S. Rubinstein and W. Hartzog, "Anonymization and Risk," *WASHINGTON LAW REVIEW*, vol. 91, p. 59, 2016.
- [7] S. Stalla-Bourdillon and A. Knight, "Anonymous data v. personal data— a false debate: An eu perspective on anonymization, pseudonymization and personal data," vol. 34, no. 2, p. 39, 2017.
- [8] S. Kung, "A Compressive Privacy approach to Generalized Information Bottleneck and Privacy Funnel problems," *Journal of the Franklin Institute*, vol. 355, no. 4, pp. 1846–1872, Mar. 2018. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S0016003217303162>
- [9] D. O'Shaughnessy, "Speaker Recognition," in *Speech Communications: Human and Machine*. IEEE, 2000, pp. 437–459. [Online]. Available: <https://ieeexplore-ieee-org.proxy.uba.uva.nl:2443/document/5312377>
- [10] N. Almaadeed, A. Aggoun, and A. Amira, "Text-Independent Speaker Identification Using Vowel Formants," *Journal of Signal Processing Systems*, vol. 82, no. 3, pp. 345–356, Mar. 2016. [Online]. Available: <http://link.springer.com/10.1007/s11265-015-1005-5>
- [11] R. J. J. H. van Son. (2019) Pseudonymize speech. Netherlands Cancer Institute. [Online]. Available: <https://robvanson.github.io/PseudonymizeSpeech/>
- [12] E. Moulines and F. Charpentier, "Pitch-synchronous waveform processing techniques for text-to-speech synthesis using diphones," *Speech communication*, vol. 9, no. 5-6, pp. 453–467, 1990.
- [13] A. C. Lammert and S. S. Narayanan, "On Short-Time Estimation of Vocal Tract Length from Formant Frequencies," *PLOS ONE*, vol. 10, no. 7, p. e0132193, 2015.
- [14] P. Boersma and D. Weenink, *Praat: a system for doing phonetics with the computer*, 2017. [Online]. Available: <http://www.praat.org>
- [15] R. J. J. H. van Son, C. Middag, and K. Demuynck, "Vowel space as a tool to evaluate articulation problems," in *Proceedings of Interspeech 2018, Hyderabad*, 2018, pp. 357–361.
- [16] N. H. De Jong and T. Wempe, "Praat script to detect syllable nuclei and measure speech rate automatically," *Behavior research methods*, vol. 41, no. 2, pp. 385–390, 2009.
- [17] R. J. J. H. van Son. (2020) Listening experiments from Adjustable Deterministic Pseudonymization of Speech. Netherlands Cancer Institute. [Online]. Available: [https://robvanson.000webhostapp.com/R\\_Van\\_Son\\_IS2020.Experiments/index.html](https://robvanson.000webhostapp.com/R_Van_Son_IS2020.Experiments/index.html)
- [18] ——. (2019) akouste. Netherlands Cancer Institute. [Online]. Available: <https://robvanson.github.io/akouste/>
- [19] R. J. J. H. Van Son, D. Binnenpoorte, H. v. d. Heuvel, and L. Pols, "The IFA corpus: a phonemically segmented Dutch "open source" speech database," in *Proceedings of EUROSPEECH 2001 Aalborg*, Aalborg, Denmark, 2001, pp. 2051–2054.
- [20] M. S. Ribeiro. (2018) Parallel audiobook corpus. [dataset]. University of Edinburgh. School of Informatics. <https://doi.org/10.7488/ds/2468>. [Online]. Available: <https://datashare.is.ed.ac.uk/handle/10283/3217>
- [21] J. Yamagishi, M. Todisco, M. Sahidullah, H. Delgado, X. Wang, N. Evans, T. Kinnunen, K. A. Lee, V. Vestman, A. Nautsch et al. (2019) Asvspoof 2019: The 3rd automatic speaker verification spoofing and countermeasures challenge database. [dataset]. University of Edinburgh. School of Informatics. <https://doi.org/10.7488/ds/2555>. [Online]. Available: <https://datashare.is.ed.ac.uk/handle/10283/3336>
- [22] R Core Team, *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria, 2019. [Online]. Available: <http://www.R-project.org/>
- [23] K. Fradette, H. J. Keselman, L. Lix, J. Algina, and R. R. Wilcox, "Conventional And Robust Paired And Independent-Samples t Tests: Type I Error And Power Rates," *Journal of Modern Applied Statistical Methods*, vol. 2, no. 2, pp. 481–496, Nov. 2003. [Online]. Available: <http://digitalcommons.wayne.edu/jmasm/vol2/iss2/22>