



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

## User Identity Protection in Automatic Emotion Recognition through Disguised Speech

**Citation for published version:**

Haider, F, Albert, P & Luz, S 2021, 'User Identity Protection in Automatic Emotion Recognition through Disguised Speech', *AI*, vol. 2, no. 4, pp. 636-649. <https://doi.org/10.3390/ai2040038>

**Digital Object Identifier (DOI):**

[10.3390/ai2040038](https://doi.org/10.3390/ai2040038)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

AI

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# User identity protection in automatic emotion recognition through disguised speech

Fasih Haider <sup>\*</sup> , Pierre Albert and Saturnino Luz 

Usher Institute, Edinburgh Medical School, The University of Edinburgh, Edinburgh EH16 4UX, UK; {Fasih.Haider, Pierre.Albert, S.Luz}@ed.ac.uk

\* Correspondence: Fasih.Haider@ed.ac.uk

**Abstract:** Ambient Assisted Living (AAL) technologies are being developed which could assist elderly people to live healthy and active lives. These technologies have been used to monitor people's daily exercises, consumption of calories and sleep patterns, and to provide coaching interventions to foster positive behaviour. Speech and audio processing can be used to complement such AAL technologies to inform interventions for healthy ageing by analyzing speech data captured in the user's home. However, the collection of data in home settings presents acute privacy protection challenges. To address this issue, we propose a low cost system for recording disguised speech signals which can protect user identity by using pitch shifting. The disguised speech so recorded can then be used for training machine learning models for affective behaviour monitoring. Affective behaviour could provide an indicator of the onset of mental health issues such as depression and cognitive impairment, and help develop clinical tools for automatically detecting and monitoring disease progression. In this article, acoustic features extracted from the non-disguised and disguised speech are evaluated in an affect recognition task using six different machine learning classification methods. The results of transfer learning from non-disguised to disguised speech are also demonstrated. We have identified sets of acoustic features which are not affected by the pitch shifting algorithm and also evaluated them in affect recognition. We found that while the non-disguised speech signal gives the best unweighted average recall (UAR) of 80.01% the disguised speech signal only causes a slight degradation in performance, reaching 76.29% UAR. The transfer learning from non-disguised to disguised speech results in a greater drop in UAR (65.13%). However, feature selection improves the UAR (68.32%). This work forms part of a large project which includes health and wellbeing monitoring and coaching.

**Keywords:** privacy preservation; affect recognition; health technologies; emotion recognition; ambient assisted living; social signal processing

**Citation:** . *Journal Not Specified* 2021, 1, 0. <https://doi.org/>

Received:

Accepted:

Published:

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Copyright:** © 2021 by the authors. Submitted to *Journal Not Specified* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Health and wellbeing monitoring using Ambient Assisted Living (AAL) technologies involves developing systems for automatically detecting and tracking a number of events that might require attention or coaching. In the SAAM project [1], we are employing AAL technologies to analyse activities and health status of older people living on their own or in assisted care settings, and to provide them with personalised multimodal coaching. Such activities and status include mobility, sleep, social activity, air quality, cardiovascular health, diet [2], emotions [3] and cognitive status [4]. While most of these signals are tracked through specialized hardware, audio and speech are ubiquitous sources of data which could also be explored in these contexts. Speech quality and activity, in particular, closely reflect health and wellbeing. We have explored the potential of speech analysis for automatically recognizing emotions [3], cognitive difficulties [4] and eating-related events [2] in the SAAM AAL environment [5]. AAL technologies and coaching systems such as SAAM, which focus on monitoring of everyday activities, can benefit from recognition of these audio events in characterizing contextual information

39 against which other monitoring signals can be interpreted. However user privacy re-  
40 mains one of the major challenges in collecting audio data in home environments for the  
41 development of health monitoring technology.

### 42 1.1. Mental Health and Affective Speech

43 The literature suggests that older people with cognitive impairment have difficulty  
44 accessing semantic information [6]. Since successful communication is essential for  
45 meaningful social interaction, this takes a toll on the patients' and their carers' wellbeing.  
46 This has an impact on the emotional life of these people. Speech monitoring for mood and  
47 cognitive changes may help inform interventions targeted at alleviating such impacts.

48 In addition to their role in cognition [7], the expression of emotions and their recog-  
49 nition are key aspects of communication [8]. Emotional information can be conveyed  
50 in different ways, from explicit facial and verbal expression (e.g. smile, pout, happy  
51 statement) to more subtle non-verbal cues, such as intonation, modulation of vocal pitch  
52 and loudness of emotional expression. These non-verbal cues are generally referred to  
53 as emotional prosody.

54 In a previous study [9], we found that there are differences in automatically inferred  
55 affective behaviours regarding expressions of *sadness*, *anger* and *disgust* among people  
56 with and without cognitive impairment (Alzheimer's Disease, AD). Although these  
57 results need further study, they suggest that speakers with AD exhibit a deficit in the  
58 expression of those emotions, reflected on voice volume, speech rate and pitch. The  
59 proposed Affective Behaviour Representation (ABR) and emotion classification scores are  
60 able to predict cognitive deficit in such situations with an accuracy of 63.42%. However,  
61 in that study there was a mismatch between the dataset used to generate the features for  
62 recognition (*emoDB* [10]) and the data on which these features were used (*Pitt Corpus*  
63 [11]). Thus, prediction accuracy is likely to have been hindered by the facts that (1) the  
64 Pitt Corpus was not explicitly designed to elicit emotions, (2) that the two datasets were  
65 recorded under different acoustic conditions, (3) that the speakers were selected from  
66 different demographics, and (4) that they are in different languages [9].

### 67 1.2. Privacy-Concerns Related to Speech

68 Privacy concerns constitute a major obstacle in developing and deploying digital  
69 technologies for monitoring cognitive health. Individual and societal concerns about  
70 privacy and data security have been translated in regulations. In the European Union,  
71 the GDPR [12] has set new standards for the collection and management of personal  
72 information. Speech data is classified as personal data<sup>1</sup>: it can be used to identify age,  
73 gender, subject identity and health status [13]. Sensitive data also encompass additional  
74 data such as content-free features which could potentially be used for the identification  
75 of a person. The potential of such features as biometric markers further widens the  
76 importance of their protection. Concern about privacy is shared by users, who are  
77 reluctant to consent to being constantly recorded at their homes and/or while speaking  
78 through phones or computers. The balance between the benefits from an analysis of  
79 spoken interaction is often offset by the associated threat to privacy.

80 Ethical requirements for health-studies have reflected these changes in regulation.  
81 They have raised awareness on the need for careful risk analysis for studies involving  
82 the collection and use of speech-related data. In the context of AAL and in-situ studies,  
83 speech analysis usually requires sending data over networks with different levels of secu-  
84 rity and associated risks, setting the additional possibility of a data breach if intercepted  
85 and compromised. While the security of the network can be improved by reducing the  
86 transit and exposure of sensitive data through a local pre-processing [14,15], the risk  
87 posed by the presence of sensitive data remains.

<sup>1</sup> As defined in Art. 4(14) of the GDPR and Article 3(13) Directive 2016/680

88 A possible way to mitigate these problems is to obfuscate the identity of the user  
89 while the data is collected by changing the pitch of their speech [16]. However, changing  
90 the signal can also degrade its analysis: pitch shifting disturbs the acoustic patterns of  
91 speech which could be indicative of cognitive impairment.

92 Hence developing a digital technology using acoustic information should take  
93 these issues into account. In this study, we also propose a framework using feature  
94 engineering to address the disturbance of acoustic features caused by pitch alteration for  
95 affect recognition as shown in Figure 5d.

### 96 1.3. Speech Disguising

97 Speech Disguising is a way to alter speech to hide someone's identity [16]. Zheng et  
98 al. [17] subjectively analyse the automatic speech disguise technologies i.e. pitch shifting,  
99 vocal tract length normalization (VTLN) and voice conversion (VC) using 30 trials. They  
100 found that the speech disguise technologies greatly confuse human evaluators, with an  
101 equal error rate around random guess (i.e. 50.00 % for pitch shifting, 46.67% for VTLN  
102 and 46.67 % for VC).

### 103 1.4. Contribution

104 We have previously developed a low-cost system [15,18] which records content-free,  
105 anonymised audio features for automatic analysis. In particular, we extract features such  
106 as the *eGeMAPS* set [19] which we have used to detect specific behaviours in the above-  
107 mentioned applications [2–4]. However, one of the limitation was that the previous  
108 system [15] delete the audio file after extracting the acoustic features from user's speech.  
109 It could work if the emotion is self-reported by user, and we do not have a plan to  
110 evaluate the new features (i.e. going to be proposed in future), but not for situations  
111 where other humans needs to annotate the audio files with emotions to generate data for  
112 machine learning model training. So that preserving audio file is also important while  
113 preserving privacy. While speech disguising technologies could help preserve the user's  
114 privacy to some extent, a question arises: "what are the effects of speech disguising on  
115 acoustic information for emotion recognition"? In this study, we extend our previous  
116 work and propose to collect the disguised speech by altering the pitch of the speech  
117 signal to protect the identity of a user for development and deployment of machine  
118 learning based application. For testing (i.e. deployment), this approach also guarantees  
119 the user's spoken content privacy in addition to identity protection. This is because the  
120 acoustic features are computed using different statistical functionals at the utterance  
121 level rather than at frame level, which makes it impossible to extract or re-build content  
122 information through, for instance, synthesis of speech from the extracted features or  
123 automatic transcription [20].

124 To the authors' best knowledge, this is the first study and evaluation of disguised  
125 speech for the development and deployment of affect recognition technologies based on  
126 acoustic features. Hence, the contributions of this article are:

- 127 • Identification of acoustic features which are not affected by disguising speech;
- 128 • Evaluation of acoustic features extracted from the disguised speech for affect recog-  
129 nition, and comparison with features extracted from non-disguised speech;
- 130 • Demonstration of transfer-learning of acoustic features from non-disguised speech  
131 to disguised speech for affect recognition, and analysis of their generalisability.

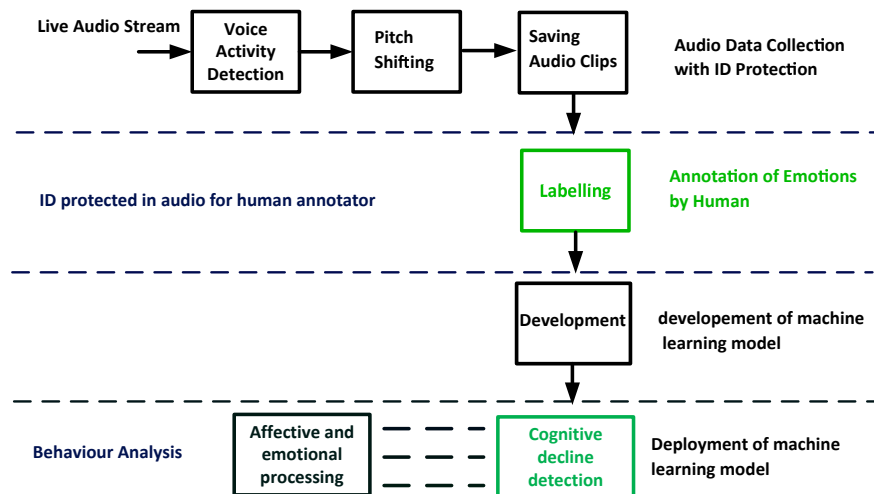
## 132 2. Materials and Methods

133 This section describes the system and algorithms which have been used for propos-  
134 ing emotion recognition using disguised speech.

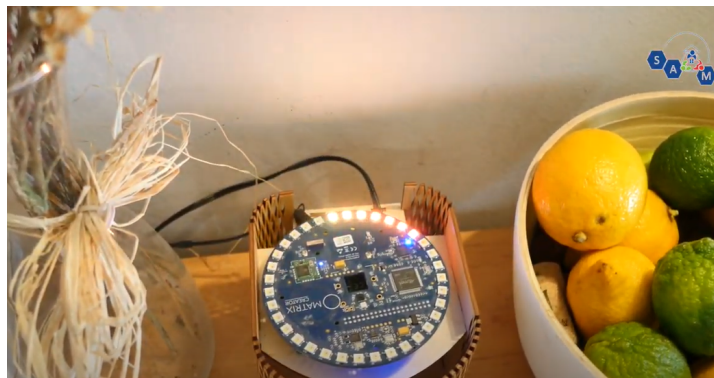
### 135 2.1. Emotion Recognition System

136 This section describes hardware and software components of the system used to  
137 extract acoustic features and collect disguised speech. The collected disguised speech

138 could be presented to human annotators (e.g. crowd-sourced annotation i.e labelling  
 139 stage) for annotation of emotions. The system's architecture is shown in Figure 1 where  
 140 the voice activity detection module detects audio segments based on energy of audio  
 141 signal. After that, we use pitch shifting algorithm [21] for speech disguising and saves  
 142 the audio segments. Later, we extract acoustic features using openSMILE [22] and train  
 143 machine learning models (development module) for emotion recognition. At the end,  
 144 we test the machine learning model (affective and emotional processing module).



**Figure 1.** Proposed approach: the affective and emotional processing module will provide input to the cognitive decline recognition module. The 'labelling' and 'cognitive decline detection' are not part of this study. The pitch shifting parameters are only known to and set by the data collection technician and/or user. The Human annotator doesn't have that information.



**Figure 2.** Matrix Creator and Raspberry Pi 3 B+

### 145 2.1.1. Hardware Components

146 The hardware consists of a Matrix Creator board, constituted of a microphone array,  
 147 an inertial measurement unit, and several other sensors, mounted on a Raspberry Pi  
 148 3 B+, as shown in Figure 2. This setup is meant to be installed in a room where social  
 149 activity and dialogue interaction occurs frequently, such as a dining room or a sitting  
 150 room.

### 151 2.1.2. Software Components

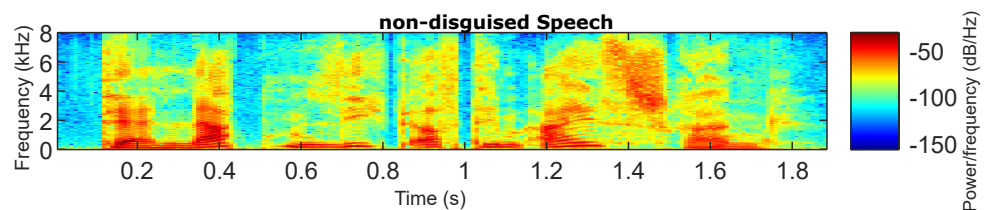
152 For voice activity detection, we employed the Auditok<sup>2</sup> Python binding. As changes  
153 are detected on disk due voice recording (using the watchdog library<sup>3</sup>) the OpenSMILE  
154 [23] toolkit processes the audio file of disguised speech and saves the speech features  
155 in the attribute-relation file format (ARFF). The extracted acoustic features are then  
156 processed by a machine learning model for emotion recognition.

### 157 2.2. Data sets

158 The Berlin Database of Emotional Speech (EmoDB) corpus [10] is a data set com-  
159 monly used in the automatic emotion recognition literature. It features 535 acted emo-  
160 tions in German (5 male and 5 females), based on utterances carrying no emotional  
161 bias. The corpus was recorded in a controlled environment resulting in high quality  
162 recordings. Actors were allowed to move freely around the microphones, which affected  
163 absolute signal intensity. In addition to the emotion, each recording was labelled with  
164 phonetic transcription using the SAMPA phonetic alphabet, emotional characteristics  
165 of the voice, segmentation of the syllables, and stress. The quality of the data set was  
166 evaluated by perception tests carried out by 20 human participants. In a first recognition  
167 test, subjects listened to a recording once before assigning one of the available categories,  
168 achieving an average recognition rate of 86%. A second naturalness test was performed.  
169 Documents achieving a recognition rate lower than 80% or a naturalness rate lower than  
170 60% were discarded from the main corpus, reducing the corpus to 535 recordings from  
171 the original 800. The data sets is annotated for 6+1 emotions: anger, disgust, fear, joy  
172 (happiness), sadness, boredom + neutral.

### 173 2.3. Identity Protection

174 To disguise the identify of the subjects, we apply pitch shifting algorithm while  
175 maintaining the duration of speech signal using Praat [21]. The audio data with iden-  
176 tity protection along with script for pitch shifting is made available through our git  
177 repository<sup>4</sup>. We have used a factor of 2 for pitch shifting with time step of 0.01 seconds,  
178 minimum pitch of 75 Hz, and maximum pitch of 600 Hz. The pitch shifting parameters  
179 are only known to and set by the data collection technician and/or user. The Human  
180 annotator does not have that information. An example of non-disguised and disguised  
181 audio segment (i.e. spectrogram representation) is shown in Figure 3 and 4 respectively,  
182 where the durations of the non-disguised and disguised speech are the same.



**Figure 3.** An example of a speech utterance’s spectrogram from the EmoDB dataset of a male subject.

### 183 2.4. Acoustic Features

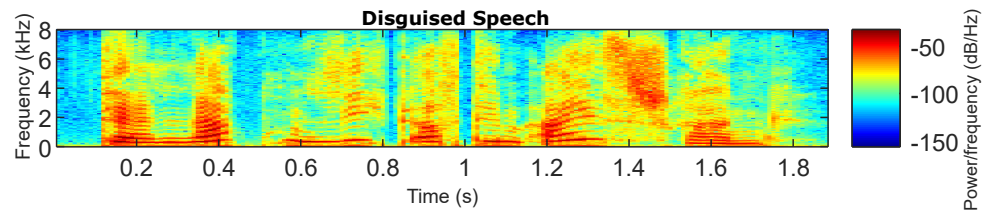
184 Acoustic feature extraction was performed on the non-disguised and disguised  
185 speech segments using the openSMILE v2.1 toolkit which is a “source-available” software  
186 suite for automatic extraction of features from speech, widely used for emotion and  
187 affect recognition in speech [24]. The extracted features are also made available through

<sup>2</sup> <https://pypi.org/project/auditok/> – accessed April 2021

<sup>3</sup> <https://github.com/gorakhargosh/watchdog> – accessed April 2019

<sup>4</sup> <https://git.ecdf.ed.ac.uk/fhaider/pitchshifting4affectrecognition>





**Figure 4.** An example of a speech utterance's spectrogram from the EmoDB dataset of a male subject after applying pitch shifting algorithm for identity protection.

188 the above mentioned git repository. The following is a brief description of the acoustic  
189 feature sets used in the experiments described in this paper:

#### 190 2.4.1. emobase:

191 This feature set contains the mel-frequency cepstral coefficients (MFCC), voice  
192 quality, fundamental frequency (F0), F0 envelope, line spectral pairs (LSP) and intensity  
193 features with their first and second order derivatives. Several statistical functions are  
194 applied to these features, resulting in a total of 988 features for every speech segment  
195 [24].

#### 196 2.4.2. ComParE:

197 The *ComParE 2013* [23] feature set includes energy, spectral, MFCC, and voicing  
198 related low-level descriptors (LLDs). LLDs include logarithmic harmonic-to-noise ratio,  
199 voice quality features, Viterbi smoothing for F0, spectral harmonicity and psychoacoustic  
200 spectral sharpness. Statistical functionals are also computed, bringing the total to 6,373  
201 features.

#### 202 2.4.3. eGeMAPS:

203 The *eGeMAPS* [19] feature set resulted from an attempt to reduce the somewhat  
204 unwieldy feature sets above to a reduced set of acoustic features based on their potential  
205 to detect physiological changes in voice production, as well as theoretical significance  
206 and proven usefulness in previous studies [22]. It contains the F0 semitone, loudness,  
207 spectral flux, MFCC, jitter, shimmer, F1, F2, F3, alpha ratio, Hammarberg index and  
208 slope V0 features, as well as their most common statistical functionals, for a total of 88  
209 features per speech segment.

### 210 2.5. Statistical Analysis

211 To investigate the possible differences in acoustic characteristics between the non-  
212 disguised and disguised speech signals, we first performed a normality test using the  
213 one-sample Kolmogorov-Smirnov procedure. This test showed that the data (i.e. acoustic  
214 features) follow a normal distribution ( $p < 0.001$ ). We then performed a t-test between  
215 the acoustic features extracted from the non-disguised speech signals and the acoustic  
216 features extracted from the disguised speech signal. We observed the following:

- 217 1. for the emobase feature set, there are 257 features out of 988 for which no statistically  
218 significant differences ( $p > 0.05$ ) between the non-disguised and disguised speech  
219 signals were found. Parts of different functional of Mfcc, fftMag, ZCR, energy,  
220 loudness and intensity are not affected by the speech alteration.
- 221 2. For the ComParE feature set, we found that 2491 features out of 6373 show no  
222 statistically significant differences ( $p > 0.05$ ) between non-disguised and disguised  
223 speech signals. Some mfcc, fftMag, audiospec, HNR, ZCR, energy, RASTA, jitter  
224 and shimmer functionals are not affected by the speech alteration procedure. The  
225 full lists of emobase and ComParE features tested is available through the above  
226 mentioned git repository.

227 3. For the eGeMAPS feature set, we have noted that there are 24 features out of 88  
 228 which have no statistically significant differences ( $p > 0.05$ ). The full list of those  
 229 features is shown below:

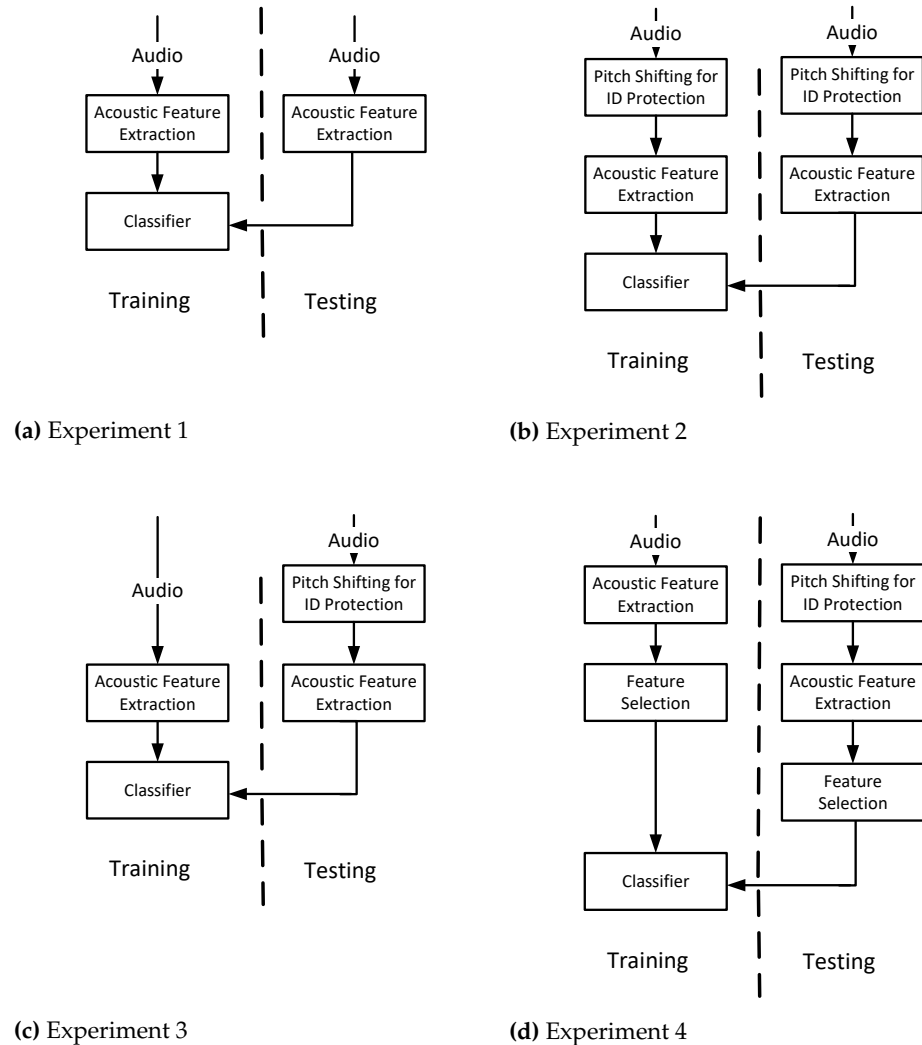
- 230 • *F0semitoneFrom27.5Hz\_sma3nz\_pctlrangle0 - 2*
- 231 • *F0semitoneFrom27.5Hz\_sma3nz\_meanRisingSlope*
- 232 • *F0semitoneFrom27.5Hz\_sma3nz\_stddevRisingSlope*
- 233 • *F0semitoneFrom27.5Hz\_sma3nz\_stddevFallingSlope*
- 234 • *loudness\_sma3\_meanRisingSlope*
- 235 • *spectralFlux\_sma3\_stddevNorm*
- 236 • *mfcc1\_sma3\_stddevNorm*
- 237 • *mfcc2\_sma3\_stddevNorm*
- 238 • *mfcc3\_sma3\_stddevNorm*
- 239 • *logRelF0 - H1 - H2\_sma3nz\_stddevNorm*
- 240 • *logRelF0 - H1 - A3\_sma3nz\_stddevNorm*
- 241 • *alphaRatioV\_sma3nz\_amean*
- 242 • *alphaRatioV\_sma3nz\_stddevNorm*
- 243 • *hammarbergIndexV\_sma3nz\_amean*
- 244 • *slopeV0 - 500\_sma3nz\_stddevNorm*
- 245 • *slopeV500 - 1500\_sma3nz\_stddevNorm*
- 246 • *spectralFluxV\_sma3nz\_stddevNorm*
- 247 • *mfcc1V\_sma3nz\_stddevNorm*
- 248 • *mfcc2V\_sma3nz\_stddevNorm*
- 249 • *mfcc3V\_sma3nz\_stddevNorm*
- 250 • *mfcc4V\_sma3nz\_stddevNorm*
- 251 • *loudnessPeaksPerSec*
- 252 • *MeanUnvoicedSegmentLength*
- 253 • *StddevUnvoicedSegmentLength*

## 254 2.6. Classification Methods

255 The classification experiments were performed using six different methods, namely  
 256 decision trees (DT, where the leaf size is optimized through a grid search within a range  
 257 of 1 to 20), nearest neighbour (KNN, where K parameter is optimized through a grid  
 258 search within a range of 1 to 10), linear discriminant analysis (LDA), random forest (RF,  
 259 with 1500 trees, where leaf size is optimized through a grid search within a range of 1 to  
 260 20), Naive Bayes (NB, with kernel distribution assumption optimized through a grid  
 261 search for kernel smoothing density estimate, Multinomial distribution, Multivariate  
 262 multinomial distribution and Normal distribution) and support vector machines: SVM,  
 263 with a linear kernel (optimized by trying different kernel function i.e., linear, Gaussian,  
 264 RBF and polynomial) with box constraint optimized by trying a grid search between  
 265 0.1 to 1.0, and sequential minimal optimization solver (optimized by trying different  
 266 solvers i.e., iterative single data algorithm, L1 soft-margin minimization by quadratic  
 267 programming and sequential minimal optimization ). The prior-probabilities of the  
 268 classifiers are set according to the class distributions.

269 The classification methods are implemented in MATLAB (<http://uk.mathworks.com/products/matlab/> (December 2020)) using the statistics and machine-learning  
 270 toolbox. The classifier hyper-parameters maximum ranges (such as  $K = 10$ ) are set  
 271 through trial and error. A leave-one-subject-out (LOSO) cross-validation setting was  
 272 adopted, where the training data does not contain any information of the validation  
 273 subjects. To assess the classification results, we used the Unweighted Average Recall  
 274 (UAR) instead of overall accuracy as the dataset is imbalanced. The unweighted average  
 275 recall is the arithmetic mean of recall for all seven classes.  
 276





**Figure 5.** Affect recognition system: Machine learning model training and testing where testing is performed in leave one subject out cross-validation settings.

### 277 3. Experimentation

278 This section describes the experiments and data partition to evaluate the proposed  
279 frameworks as shown in Figure 5.

#### 280 3.1. Experiment 1

281 In this experiment, we extracted acoustic features over the non-disguised audio  
282 data. Later we trained the machine learning models for classification purpose. The  
283 validation is performed in leave-one subject out cross-validation setting as shown in  
284 Figure 5a.

#### 285 3.2. Experiment 2

286 In this experiment, we extracted acoustic features over the transformed audio data  
287 where we hid the identity of a subject using pitch shifting algorithm. Later we trained  
288 the machine learning models for classification purpose. The validation is performed in  
289 leave-one subject out cross-validation setting as shown in Figure 5b.

### 290 3.3. Experiment 3

291 In this experiment, we trained the machine learning models using non-disguised  
292 speech and the validation is performed using disguised speech in leave-one subject out  
293 cross-validation setting as shown in Figure 5c.

### 294 3.4. Experiment 4

295 This experiment uses the selected acoustic features as described in Section 2.5, we  
296 trained the machine learning models using non-disguised speech and the validation is  
297 performed using disguised speech in leave-one subject out cross-validation setting as  
298 shown in Figure 5d.

## 299 4. Results

300 This section reports the results for the four experiments.

### 301 4.1. Experiment 1

302 The UAR for all feature sets and classification methods is shown in Tables 1. These  
303 results indicate that the ComParE feature set (80.01%) provides the best UAR, with the  
304 LDA classifier for emotion recognition. The confusion matrix is shown in Figure 6 for  
305 further insight (i.e. precision and recall for all 6+1 emotions) into the best result. The  
306 results indicate that the SVM provides the best averaged UAR of 73.42% across all the  
307 feature sets, and the ComParE feature set (57.76%) provides the best average UAR across  
the all classifiers.

Table 1: Experiment 1: Affect recognition results without identity protection where training and validation is performed on the non-disguised audio data. The Unweighted Average Recall (UAR%) is reported.

Features	RF	DT	KNN	NB	SVM	LDA	avg.
emobase.	0.6835	0.5052	0.2460	0.6051	0.7308	0.5574	0.5547
ComParE	0.7059	0.5368	0.2281	0.3953	0.7949	<b>0.8001</b>	<b>0.5768</b>
eGeMAPS	0.7063	0.4918	0.3885	0.4854	0.6858	0.6616	0.5699
avg	0.6986	0.5113	0.2875	0.4953	<b>0.7372</b>	0.6730	—

308

True Class	Predicted Class							Recall
	Anger	Bore.	Disgust	Fear	Happy	Sad	Neutral	
Anger	115		1		11			90.6%
Bore.		72				2	7	88.9%
Disgust			38	3	2	1	2	82.6%
Fear	5		1	49	9	1	4	71.0%
Happy	23		2	5	41			57.7%
Sad		4	1	2		50	5	80.6%
Neutral		5	2	1		1	70	88.6%
Precision	80.4%	88.9%	84.4%	81.7%	65.1%	90.9%	79.5%	

UAR = 80.01%  
Accuracy = 81.31%

Figure 6. Confusion matrix of the best result for experiment 1 using LDA and Compare Feature set.

### 309 4.2. Experiment 2

310 The UAR for all feature sets and classification methods is shown in Table 2. These  
311 results indicate that the combination of the ComParE feature set and LDA again provides  
312 the best UAR score (76.29%). The confusion matrix for this is shown in Figure 7 where  
313 precision and recall for all 6+1 emotions are listed. In addition, SVM provides the best

314 averaged UAR of 71.68% across all the feature sets and the eGeMAPS feature set (54.78%)  
 315 provides the best average UAR across the all classifiers.

Table 2: Experiment 2: Affect recognition results with identity protection for training and validation subjects where training and validation is performed on the pitch-shifted audio data. The Unweighted Average Recall (UAR%) is reported.

Features	RF	DT	KNN	NB	SVM	LDA	avg.
emobase.	0.6657	0.4588	0.2759	0.5865	0.7358	0.5417	0.5441
ComParE	0.7063	0.5211	0.2016	0.2440	0.7388	<b>0.7629</b>	0.5291
eGeMAPS	0.6335	0.4529	0.3705	0.4818	0.6759	0.6720	<b>0.5478</b>
avg	0.6685	0.4776	0.2827	0.4374	<b>0.7168</b>	0.6589	—

True Class	Predicted Class							Recall
	Anger	Bore.	Disgust	Fear	Happy	Sad	Neutral	
Anger	103		1	4	19			81.1%
Bore.	1	68		1	1	1	9	84.0%
Disgust	1	1	35	1	1	1	6	76.1%
Fear	9		2	48	4	2	4	69.6%
Happy	26		1	6	37		1	52.1%
Sad		1	1	4		52	4	83.9%
Neutral		6	3		1		69	87.3%
	Precision	73.6%	89.5%	81.4%	75.0%	58.7%	92.9%	74.2%

UAR = 76.29%  
Accuracy = 77.01%

Figure 7. Confusion matrix of the best result for experiment 2 using LDA and Compare Feature set.

#### 316 4.3. Experiment 3

317 The results for this experiment are shown in Tables 3. These results indicate that  
 318 the ComParE feature set again provides the best UAR (65.13%), but this time the RF  
 319 classifier proves to be the most effective. The confusion matrix is shown in Figure 8  
 320 where precision and recall for all 6+1 emotions are listed. RF provides the best averaged  
 321 UAR of 57.33% across all feature sets, and the emobase feature set yields the best average  
 322 UAR across all classifiers (45.95%).

Table 3: Experiment 3: Affect recognition results with identity protection for validation subjects, where training is performed on the non-disguised audio data and validation is performed on the pitch-shifted audio data. The Unweighted Average Recall (UAR%) is reported.

Features	RF	DT	KNN	NB	SVM	LDA	avg.
emobase.	0.5624	0.4172	0.2162	0.4838	0.6103	0.4673	<b>0.4595</b>
ComParE	<b>0.6513</b>	0.4479	0.2161	0.1429	0.1435	0.1344	0.2893
eGeMAPS	0.5062	0.3698	0.2623	0.3470	0.5391	0.1339	0.3597
avg	<b>0.5733</b>	0.4116	0.2315	0.3246	0.4310	0.2452	—

#### 323 4.4. Experiment 4

324 The resulting UAR scores for all feature sets and classification methods used in this  
 325 experiment are shown in Tables 4. As before, the ComParE/RF combination achieves  
 326 the best result (68.32%). The confusion matrix is shown in Figure 9 where precision  
 327 and recall for all 6+1 emotions are listed. As in the previous experiment, RF provided  
 328 the best averaged UAR (60.34%) across all the feature sets, and the emobase feature set  
 329 yielded the best average UAR across classifiers (48.62%).

True Class	Predicted Class							Recall
	Anger	Bore.	Disgust	Fear	Happy	Sad	Neutral	
Anger	125			2				98.4%
Bore.	5	62	2	4	2	1	5	76.5%
Disgust	11	3	17	7	7	1		37.0%
Fear	19		1	40	8	1		58.0%
Happy	36		1	3	31			43.7%
Sad		6		4		49	3	79.0%
Neutral	3	8	1	5	12		50	63.3%
Precision	62.8%	78.5%	77.3%	63.5%	50.0%	94.2%	86.2%	UAR = 65.13%
								Accuracy = 69.91%

**Figure 8.** Confusion matrix of the best result for experiment 3 using RF and Compare Feature set.

Table 4: Experiment 4: Affect recognition results with identity protection, where training and validation is performed on selected acoustic features of the non-disguised audio data and validation is performed on the pitch-shifted audio data. The Unweighted Average Recall (UAR%) is reported.

Features	RF	DT	KNN	NB	SVM	LDA	avg.
emobase.	0.5731	0.4121	0.2665	0.5331	0.6250	0.5075	<b>0.4862</b>
ComParE	<b>0.6832</b>	0.4793	0.2541	0.1429	0.1839	0.1231	0.3111
eGeMAPS	0.5540	0.4467	0.2623	0.3305	0.4988	0.4375	0.4216
avg	<b>0.6034</b>	0.4460	0.2610	0.3355	0.4359	0.3560	—

True Class	Predicted Class							Recall
	Anger	Bore.	Disgust	Fear	Happy	Sad	Neutral	
Anger	125			2				98.4%
Bore.	2	70		1		1	7	86.4%
Disgust	8	7	23	7		1		50.0%
Fear	16	1	2	44	3	1	2	63.8%
Happy	44	2	2	5	13		5	18.3%
Sad		4		4		49	5	70.0%
Neutral	1	6		6	1		65	82.3%
Precision	63.8%	77.8%	85.2%	63.8%	76.5%	94.2%	77.4%	UAR = 68.32%
								Accuracy = 72.71%

**Figure 9.** Confusion matrix of the best result for experiment 4 using RF and Compare Feature set.

### 330 5. Discussion

331 The summary of results is shown in Table 5. We note that the non-disguised speech  
 332 (i.e. Experiment 1) provides the best UAR and accuracy but experiments 4 and 3 provide  
 333 the best recall for Anger (98.43%) and Sad (83.87%) as shown in bold in Table 5. The  
 334 ‘Happy’ emotion is miss-classified as ‘Anger’ and the miss-classification rate increases  
 335 for disguised speech experiments, with the worst miss-classification rate occurring when  
 336 feature selection is performed (Experiment 4). A similar patten is observed for the  
 337 ‘Disgust’ category, which exhibits the greatest performance degradation in disguised  
 338 speech. However, feature selection provides better overall UAR (68.32%) than the full  
 339 feature set (65.13%). Experiment 2 provides better UAR (76.29%) than experiments 3  
 340 and 4. One of the advantages of the architecture employed in experiment 2 is that the  
 341 training and testing are both performed on the disguised speech, with the pitch shifted  
 342 by the same factor (i.e. 2) for all speech utterances. A variable pitch factor may result in  
 343 a different outcome.

344 To better understand the relationship between the experiments, we also plotted the  
 345 Venn diagram shown in Figure 10. In this diagram, the brown area (labelled “Target”)

Table 5: Results Summary: Accuracy (Accu.), Unweighted Average Recall (UAR) and recall of each emotion for the best best results of each experiment

Experiment	Accu.	UAR	Anger	Bore.	Disgust	Fear	Happy	Sad	Neutral
EXP.1	<b>81.31</b>	<b>80.01</b>	90.55	<b>88.89</b>	<b>82.61</b>	<b>71.01</b>	<b>57.75</b>	80.65	88.61
EXP.2	77.01	76.29	81.10	83.95	76.09	69.57	52.11	<b>83.87</b>	<b>87.34</b>
EXP.3	69.91	65.13	<b>98.43</b>	76.54	36.96	57.97	43.66	79.03	63.29
EXP.4	72.71	68.32	<b>98.43</b>	86.42	50.00	63.77	18.31	79.03	82.28

346 represents the annotated labels, the blue area represents the predicted labels of *Experiment*  
 347 *1*, the red area represents the predicted labels of Experiment 2, the green area represents  
 348 the prediction obtained with the experiment 3, and finally the yellow area represents  
 349 labels predicted with the experiment 4. The Venn diagrams suggest the information  
 350 captured by different pitch profiles is not similar, as only 289 out of 535 instances are  
 351 detected by all the experiments.

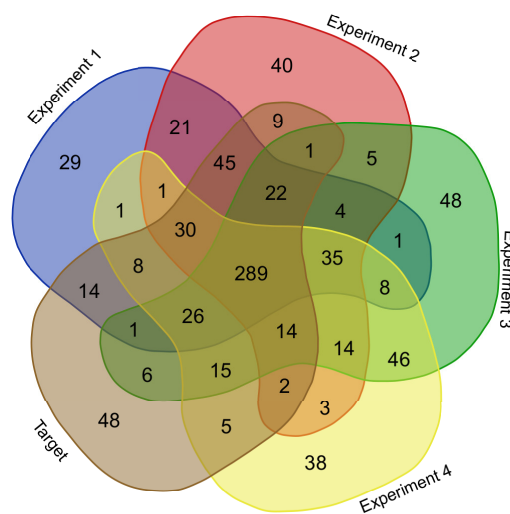


Figure 10. Venn Diagram.

352 Overall, the experiments show that despite some degradation in prediction accuracy,  
 353 privacy preservation is compatible with emotion recognition in the settings proposed.  
 354 We note that while previous studies have proposed affect recognition systems [19,25–28],  
 355 this study presents an analysis of affect recognition on data that have been transformed  
 356 to protect the identity of users.

### 357 5.1. Limitations

358 Some limitation of this study which we intend to address in future work include:

- 359 • the use of an off-the-shelf pitch shifting method which could have an influence on  
 360 the performance of affect recognition system;
- 361 • the fact that pitch is shifted using a constant factor of 2, whereas a different factor  
 362 or a variable factor could result in different results;
- 363 • feature selection is performed through a statistical approach, and more sophisticated  
 364 feature selection methods [27] might improve the results further;
- 365 • the disguised speech for affect recognition system is evaluated using data which is  
 366 collected in lab-settings instead of real-world settings;
- 367 • the hardware used for the proposed system is a combination of matrix creator and  
 368 Raspberry Pi 3 B+ with 1.4 GHz 64-bit quad-core processor, which limits one's  
 369 choice of audio processing and features extraction algorithms due to performance  
 370 limitations.

## 371 6. Conclusion

372 AAL can benefit from unobtrusive, privacy-preserving systems for gathering and  
 373 processing of speech at home. This paper described a framework for capturing disguised  
 374 speech and training machine learning models while protecting the identity of users  
 375 for automatic wellbeing monitoring tasks, in the context of an AAL-based coaching  
 376 system for healthy ageing. This study also demonstrated that the acoustic information  
 377 of disguised speech can be used for emotion recognition. We found that while the  
 378 non-disguised speech signal gives the best Unweighted Average Recall (UAR) of 80.01%  
 379 the disguised speech signal only causes a slight degradation of performance, reaching  
 380 76.29%. The transfer learning from non-disguised to disguised speech results in a  
 381 reduction of UAR (65.13%). However, feature selection improves the UAR (68.32%).  
 382 Privacy protection and preservation in audio and speech can be regarded from different  
 383 perspectives, including the protection of a person's identity, protection of the content  
 384 spoken, and protection from inferences one may be able to draw from the characteristics  
 385 of a person's voice (such as cognitive or emotional status) [29]. A current limitation of  
 386 the pitch shifting approach is that it only addresses the first (using pitch shifting for  
 387 identity protection) and second (using statistical functionals of acoustic features instead  
 388 of content) of these aspects. In future, we aim to address inference protection within  
 389 a general framework. We also plan to evaluate humans' annotation performance on  
 390 disguised speech.

391 **Author Contributions:** Conceptualization, Fasih Haider, Pierre Albert and Saturnino Luz; Data  
 392 curation, Fasih Haider and Pierre Albert; Formal analysis, Fasih Haider; Funding acquisition,  
 393 Saturnino Luz; Investigation, Fasih Haider, Pierre Albert and Saturnino Luz; Methodology, Fasih  
 394 Haider; Project administration, Saturnino Luz; Software, Fasih Haider; Supervision, Saturnino  
 395 Luz; Writing – original draft, Fasih Haider; Writing – review & editing, Fasih Haider, Pierre Albert  
 396 and Saturnino Luz.

397 **Funding:** This research has received funding from the European Union's Horizon 2020 research  
 398 and innovation programme under grant agreement No 769661, SAAM project.

399 **Conflicts of Interest:** the Authors declare no conflict of interest.

## References

1. Dimitrov, Y.; Gospodinova, Z.; Žnidaršič, M.; Ženko, B.; Veleva, V.; Miteva, N. Social Activity Modelling and Multimodal Coaching for Active Aging. *Procs. of Personalized Coaching for the Wellbeing of an Ageing Society, COACH'2019*, 2019.
2. Haider, F.; Pollak, S.; Zarogianni, E.; Luz, S. SAAMEAT: Active Feature Transformation and Selection Methods for the Recognition of User Eating Conditions. *Proceedings of the 20th ACM International Conference on Multimodal Interaction*; ACM: New York, NY, USA, 2018; ICMI '18, pp. 564–568. doi:10.1145/3242969.3243685.
3. Haider, F.; Luz, S. Attitude Recognition Using Multi-resolution Cochleagram Features. *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 3737–3741. doi:10.1109/ICASSP.2019.8682974.
4. Luz, S.; la Fuente, S.D. A Method for Analysis of Patient Speech in Dialogue for Dementia Detection. *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*; Kokkinakis, D., Ed.; European Language Resources Association (ELRA): Paris, France, 2018.
5. Hrovat, A.; Znidarsic, M.; Zenko, B.; Vucnik, M.; Mohorcic, M. Saam: Supporting active ageing-use cases and user-side architecture. *2018 27th European Conference on Networks and Communications (EuCNC)*, 2018.
6. Bondi, M.W.; Salmon, D.P.; Kaszniak, A.W. The neuropsychology of dementia. In *Neuropsychological assessment of neuropsychiatric disorders.*, 2 ed.; Oxford University Press: New York, NY, US, 1996; pp. 164–199.
7. Hart, R.P.; Kwentus, J.A.; Taylor, J.R.; Harkins, S.W. Rate of forgetting in dementia and depression. *Journal of Consulting and Clinical Psychology* **1987**, *55*, 101–105.
8. Lopes, P.N.; Brackett, M.A.; Nezlek, J.B.; Schütz, A.; Sellin, I.; Salovey, P. Emotional intelligence and social interaction. *Personality and social psychology bulletin* **2004**, *30*, 1018–1034.
9. Haider, F.; De La Fuente Garcia, S.; Albert, P.; Luz, S. Affective Speech for Alzheimer's Dementia Recognition. *LREC: Resources and Processing of linguistic, para-linguistic and extra-linguistic Data from people with various forms of cognitive/psychiatric/developmental impairments (RaPID)*; Kokkinakis, D.; Lundholm Fors, K.; Themistocleous, C.; Antonsson, M.; Eckerström, M., Eds. *European Language Resources Association (ELRA)*, 2020, pp. 67–73.
10. Burkhardt, F.; Paeschke, A.; Rolfes, M.; Sendlmeier, W.F.; Weiss, B. A database of German emotional speech. *Proceedings of the ninth European Conference on Speech Communication and Technology*, 2005, pp. 1516–1520.



11. Becker, J.; Boller, F.; Lopez, O.; Saxton, J.; McGonigle, K. The natural history of Alzheimer's disease: Description of study cohort and accuracy of diagnosis. *Archives of Neurology* **1994**, *51*, 585–594.
12. Parliament, T.E.; the Council. Regulation (EU) 2016/679 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation) **2016**.
13. Nautsch, A.; Jiménez, A.; Treiber, A.; Kolberg, J.; Jasserand, C.; Kindt, E.; Delgado, H.; Todisco, M.; Hmani, M.A.; Mtibaa, A.; others. Preserving privacy in speaker and speech characterisation. *Computer Speech & Language* **2019**, *58*, 441–480.
14. Dimitrievski, A.; Zdravevski, E.; Lameski, P.; Trajkovik, V. Addressing Privacy and Security in Connected Health with Fog Computing. Proceedings of the 5th EAI International Conference on Smart Objects and Technologies for Social Good. Association for Computing Machinery, 2019, GoodTechs '19, p. 255–260. doi:10.1145/3342428.3342654.
15. Haider, F.; Luz, S. A System for Real-Time Privacy Preserving Data Collection for Ambient Assisted Living. INTERSPEECH, 2019, pp. 2374–2375.
16. Perrot, P.; Aversano, G.; Chollet, G. Voice disguise and automatic detection: review and perspectives. *Progress in nonlinear speech processing* **2007**, pp. 101–117.
17. Zheng, L.; Li, J.; Sun, M.; Zhang, X.; Zheng, T.F. When Automatic Voice Disguise Meets Automatic Speaker Verification. *IEEE Transactions on Information Forensics and Security* **2020**, *16*, 824–837.
18. Haider, F.; Luz, S. Affect Recognition Through Scalogram and Multi-Resolution Cochleagram Features. Proc. Interspeech 2021, 2021, pp. 4478–4482. doi:10.21437/Interspeech.2021-1761.
19. Eyben, F.; Scherer, K.R.; Schuller, B.W.; Sundberg, J.; André, E.; Busso, C.; Devillers, L.Y.; Epps, J.; Laukka, P.; Narayanan, S.S.; others. The Geneva minimalistic acoustic parameter set (GeMAPS) for voice research and affective computing. *IEEE Transactions on Affective Computing* **2016**, *7*, 190–202.
20. Lajmi, L. An Improved Packet Loss Recovery of Audio Signals Based on Frequency Tracking. *Journal of the Audio Engineering Society* **2018**, *66*, 680–689.
21. Boersma, P.; Weenink, D. Praat: doing phonetics by computer [Computer program]. Version 6.0. 37. URL <http://www.praat.org/>. Retrieved March **2018**, *14*, 2018.
22. Eyben, F.; Scherer, K.R.; Schuller, B.W.; Sundberg, J.; André, E.; Busso, C.; Devillers, L.Y.; Epps, J.; Laukka, P.; Narayanan, S.S.; others. The Geneva minimalistic acoustic parameter set (GeMAPS) for voice research and affective computing. *IEEE Transactions on Affective Computing* **2016**, *7*, 190–202.
23. Eyben, F.; Wenginger, F.; Groß, F.; Schuller, B. Recent developments in opensmile, the munich open-source multimedia feature extractor. Proceedings of the 21st ACM international conference on Multimedia. ACM, Association for Computing Machinery, 2013, pp. 835–838.
24. Eyben, F.; Wöllmer, M.; Schuller, B. Opensmile: the munich versatile and fast open-source audio feature extractor. Proceedings of the 18th ACM international conference on Multimedia. ACM, Association for Computing Machinery, 2010, pp. 1459–1462.
25. Haider, F.; Luz, S. Attitude recognition using multi-resolution cochleagram features. ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2019, pp. 3737–3741.
26. Haider, F.; Salim, F.A.; Conlan, O.; Luz, S. An Active Feature Transformation Method for Attitude Recognition of Video Bloggers. INTERSPEECH, 2018, pp. 431–435.
27. Haider, F.; Pollak, S.; Albert, P.; Luz, S. Emotion recognition in low-resource settings: An evaluation of automatic feature selection methods. *Computer Speech & Language* **2020**, p. 101119.
28. Haider, F.; Pollak, S.; Albert, P.; Luz, S. Extracting Audio-Visual Features for Emotion Recognition Through Active Feature Selection. 2019 IEEE Global Conference on Signal and Information Processing (GlobalSIP). IEEE, 2019, pp. 1–5.
29. Pathak, M.A.; Raj, B.; Rane, S.D.; Smaragdis, P. Privacy-preserving speech processing: cryptographic and string-matching frameworks show promise. *IEEE signal processing magazine* **2013**, *30*, 62–74.