

An intelligent speaker that recognizes your voice

Alvin Cao, Ege Gurmericliler, Emilien Pilloud

University of Washington

Abstract

In an era where intelligent home speakers are getting increasingly popular, having an integrated personal assistant can be a decisive factor for a client looking to purchase such a device. Speaker recognition is one important technology that can create a truly personalized experience in that regard. We have implemented a system that not only recognizes commands, as many personal assistants do, but also identifies the person who speaks, as to provide a service that is tailored to the speakers very own profile.

We have designed a system that, after a short calibration period that requires users to speak, can listen to user commands and identify the registered user who speaks. For the purpose of illustration, we have implemented two commands, which are "play my playlist" and "play something I like". Each user has a dedicated playlist of songs and a list of preferred genres stored in a local database. The playlist is hardcoded, but musical genres can be added on the go with the command ["I like" + genre].

Using past research, we have gathered multiple models that have shown to be appropriate for speaker recognition. Our work distinguishes itself from past speaker recognition research in that the final system is optimized to meet constraints that emerge from real-life situations. These constraints include, and are not limited to, requiring the smallest amount of calibration samples, predicting the speaker with minimal delay and dealing with ambient noise. We have used multiple tools such as deep learning models, noise reduction, sample rate conversion and feature extraction. Our comparative analysis focuses mainly on finding the right parameters for these tools, which help us reach our above goals in the best way.

Our implementation can be found at: github.com/Emilien-P/sonos-challenge

1. Objectives

At the beginning of the project, we wanted to understand what type of challenges might occur in environments involving the Sonos speaker. Because Sonos emphasizes on having multiple connected speakers, we believe that their core demographic is groups of people who often share a common physical space, such as families and coworkers. In such scenarios, having multiple people in a space should not prevent each of them to have a personalized interaction with the Sonos speakers. We think speaker recognition can help achieve a tailored user experience for every single individual in the room. Our main goal is thus to motivate what a user can do if the Sonos speaker can recognize who he/she is.

We first started searching for tools that might help us extract features from voice samples. Past research point out the particular efficiency of Mel-Frequency Cepstrum Coefficients, which we decided to use as basis for processing our input waveforms into vectorial data (see technical definition 1) that could be used for deep learning. Our comparative analysis includes the use of additional parameters such as delta coefficients appended to MFCC coefficients (see technical definition 2). Delta coefficients give information about the MFCC trajectories. Thus, we are expecting the delta coefficients to form new vectors that include more features about the speech, and hence our models to be more accurate if trained with these new vectors.

To keep our research close to real-life scenarios, we took into consideration the ambient noise that might occur in crowded spaces. Thus, we thought about integrating a noise reduction filter in our system. Common sense made us give this immediate decision, but we soon realized that many parameters needed to be taken into consideration. Although our ears might enjoy noise reduction at some well-known frequencies, the training models might prefer reduction at different frequencies. They also might prefer no reduction at all and work best with the raw samples. Our comparative analysis takes into considerations those points and determines whether noise reduction is truly necessary for our system. We are expecting to require noise reduction in order to avoid our models to learn wrong patterns that might emerge from noise during the calibration session, or conversely, inaccurately predict users

36 because of additional noise that weren't present during calibration. A good
37 implementation of a lowpass filter with a proper cutoff is proved to improve
38 the acoustic performances[1].

39
40 Another well known real-life problem is the computation time of our sys-
41 tem. In order to keep the user's attention and flow of thought uninterrupted,
42 the response time of our application shouldn't exceed 10 seconds, and ideally
43 be between 0 and 2 seconds [2] . It is also a natural initiative to always
44 focus on providing the best effort on computation speed. The two main time
45 consuming parts of our system are computing the Mel-Frequency Cepstrum
46 coefficients and running the a prediction on our model. Both can take advan-
47 tage of smaller inputs. A smaller input signal would induce smaller number
48 of MFC coefficient computations. Consequently, the number of MFC Coeffi-
49 cients would be smaller, and hence the model would be trained with shorter
50 vectors. To do this, we want to implement a sample rate converter on the
51 input signal before computing the MFC Coefficients. We are expecting to
52 analyze and determine an optimal sample rate such that we minimize com-
53 putation time while not losing accuracy on predictions.

54
55 The speaker recognition challenge is not new, especially in the research
56 area. Researchers have been using machine learning for a long time to classify
57 speech features with techniques such as Support Vector Machine or Gaussian
58 Mixture Models. During the last decade, deep learning networks, including
59 long short-term memory networks, became increasingly popular. However,
60 when it comes to applying those methods to a real life product, a new chal-
61 lenge arises. Most of the research use a high number of samples whereas we
62 ideally want to train our model on as few samples as possible of a few seconds
63 each. Indeed, we want our system to be as user friendly as possible and ask-
64 ing for more than 8 to 12 calibration samples per user seems too cumbersome.

65
66 In our project, we chose four different machine learning model approaches
67 to see which one would be the most appropriate to our product. We selected
68 the linear Support Vector Classifier both because it is a reliable and base-
69 line method and that it can train on a little number of samples and still
70 achieve a decent accuracy. Then, we selected three different architectures of
71 neural networks that are shown on the figure below. The seqNN is a model
72 that does not consider the dimensionality of the Mel Frequency coefficients
73 matrix, i.e. it flattens the input before training. The convolutional neural

network is an approach well-known for its success in 2D images classifying. The LSTM network is an implementation of the cells described in [3] which as a memory implementation that avoids the major flaw of Recursive Neural Nets. (i.e. an explosion of the back-propagation gradient [cite] as the time-domain dimension increases)

To remedy the low number of samples, we use a method of sample generation which consists in sliding a fixed-sized window over the time-domain of the Mel coefficients matrix with some overlap. It selects overlapping sets of frames and hence increases the number of samples with the same amount of data.

We used the Scikit library [4] to implement the SVM model and the Keras [5] wrapper for Tensorflow to implement our deep learning models. Deep learning models being very architecture dependent, we did our best to design the best architectures possible remaining fairly simple but we are not able to know if they are optimally suiting for speaker recognition.

1.1. Technical definitions

(1) To calculate the MFC Coefficients, we used the speechpy library[6] to calculate 13 MFC Coefficients per frame for 44100 Hz input signals. We selected mostly typical parameters in the calculation, but a deliberate choice was made to set the frame stride equal to the frame length (2 ms), resulting in zero effective frame overlap, which we believed to be more appropriate for the purposes of classification.

(2) The Delta coefficients form the first-order frame-to-frame difference of the Mel-Frequency Coefficients. This provides the trajectories of the coefficients, hence appending delta coefficients to the MFCC vector should provide additional information about the dynamics of the MFCC. Delta coefficients can be found with the following equation:

$$d_t = \frac{\sum_{n=1}^N n(c_{t+n} - c_{t-n})}{2 \sum_{n=1}^N n^2} \quad (1)$$

103 2. Results

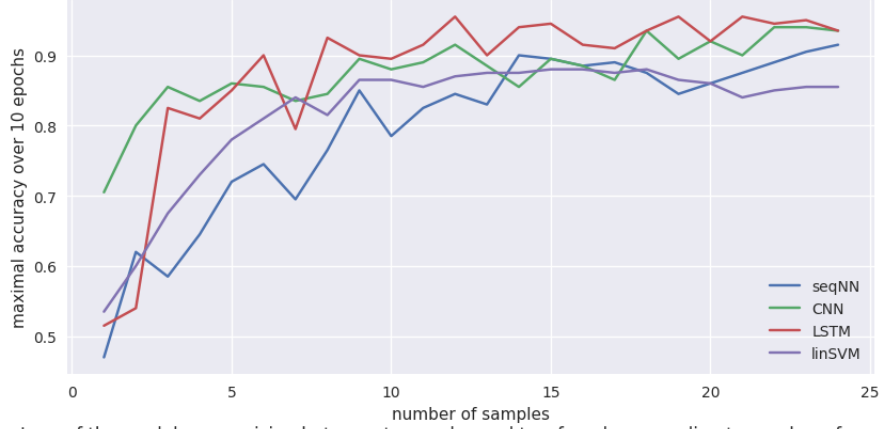
104 2.1. Machine Learning model selection

105 We implemented the four aforementioned machine learning models and
106 we got different performance rate for each of them. We used the data from
107 the CMU Festvox arctic database [7] which consist in 7 different speakers
108 male and female saying several seconds long sentences from the open-source
109 Gutenberg project. The bank contains around 500 hundred sentences of 2 to
110 5 seconds per speaker. The number of samples in the all the figures below
111 refers to the number of samples taken per user for the training. The testing
112 was made on other samples unused for training, typically around 50 testing
113 samples per speaker to get a relevant testing accuracy score.

114
115 Unsurprisingly, we could firstly see in figure 1 that the LSTM model is the
116 most performing model for our range of original sample number of interest.
117 It has above .9 accuracy level from 8 samples on with the sample extension
118 feature on. LSTM networks are well-known for their ability to classify se-
119 quences such as audio as they have memory. The seqNN architecture is the
120 overall least accurate which is also reasonable given that it loses the time di-
121 mensionality information whereas CNN still keep some 2D locality features
122 information.

123
124 In figure 2, we can see that all models do converge to a very high accuracy
125 if we train them on enough data (similar accuracy from 150 samples onwards).
126 It is a good sign that our models are not flawed.

Accuracy of the models recognizing between two males and two females according to number of samples



Loss of the models recognizing between two males and two females according to number of samples

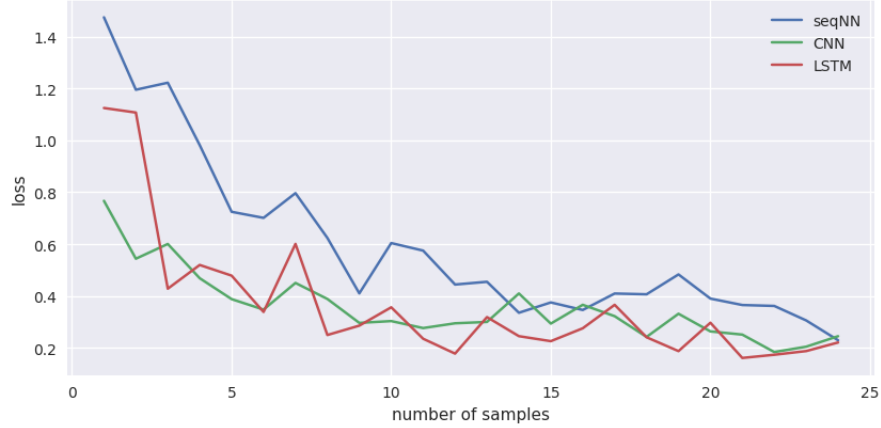


Figure 1: Accuracy of the models with sample extension

127 For small numbers of samples, we can see that the models are somewhat
 128 volatile. An important thing to notice is that both CNN and LSTM achieve
 129 significantly higher accuracy (both above 80 percent at 3 samples compared
 130 to 70 percent for linSVM and 60 percent for seqNN). All four models achieve
 131 higher accuracy as the number of samples increase, yet CNN and LSTM
 132 appear to be better even then (95percent for both compared to 90 percent
 133 and 85 percent for seqNN and linSVM respectively).

Accuracy and loss over number of samples without samples extension

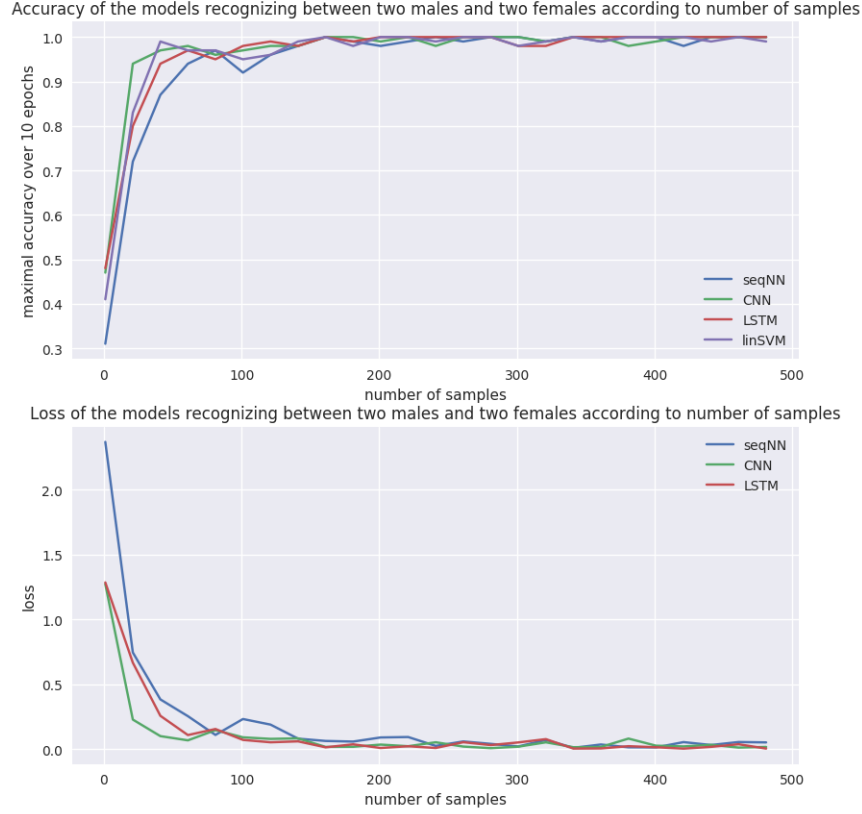


Figure 2

134 We could also conclude that our overlapping window sample extension
 135 method is useful in some cases as seen on the figure 3. The dotted lines
 136 represent the models trained on the original samples and the full lines are
 137 models also trained on the extra generated samples from windowing. The
 138 models trained on the extended data all have a higher tested accuracy around
 139 our number of samples of interest. In fact, they all start and remain at high
 140 levels of accuracy (above 75 percent). We can also see that windowing reduces
 141 volatility in small number of samples (overall between 75 and 95 percent with
 142 windowing, compared to 51 and 91 percent without).

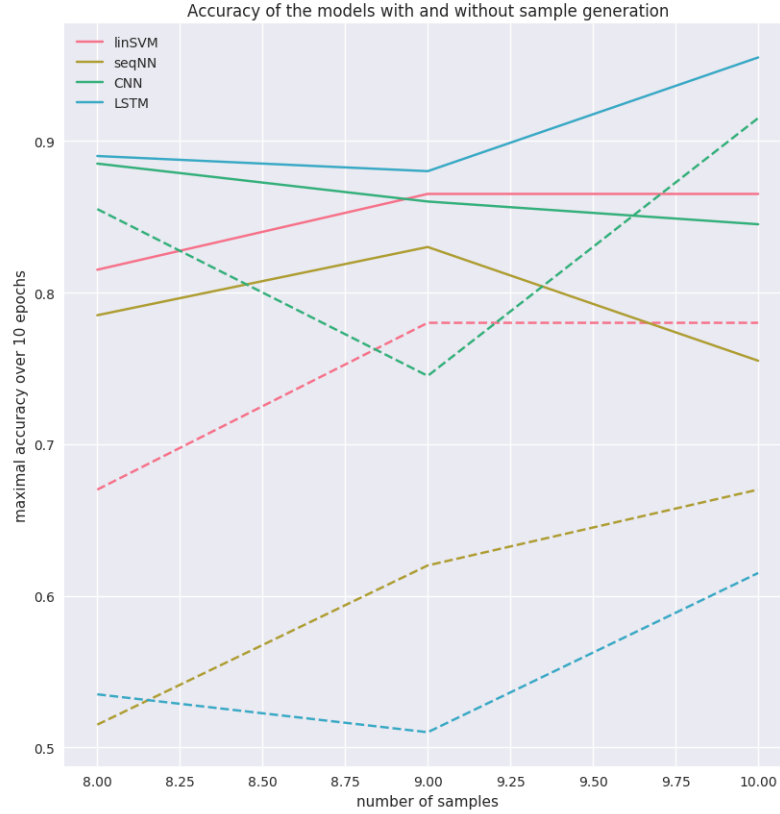


Figure 3

143 In figure 4, we plotted the accuracy of the models for different number
 144 of speakers from the arctic dataset. The experience was done three times
 145 and the bold line stands for the median accuracy result obtained. All the
 146 accuracy levels were measured for 10 training samples per speaker with the
 147 sample extension activated and 50 testing samples per speaker. We can notice
 148 once again that not only is LSTM more accurate for any number of speakers,
 149 but it also is more robust to an increase in the number of speakers. We can
 150 see the flaw of our baseline method, linear Support Machine Classifier, whose
 151 accuracy drops drastically as the number of speakers increases.

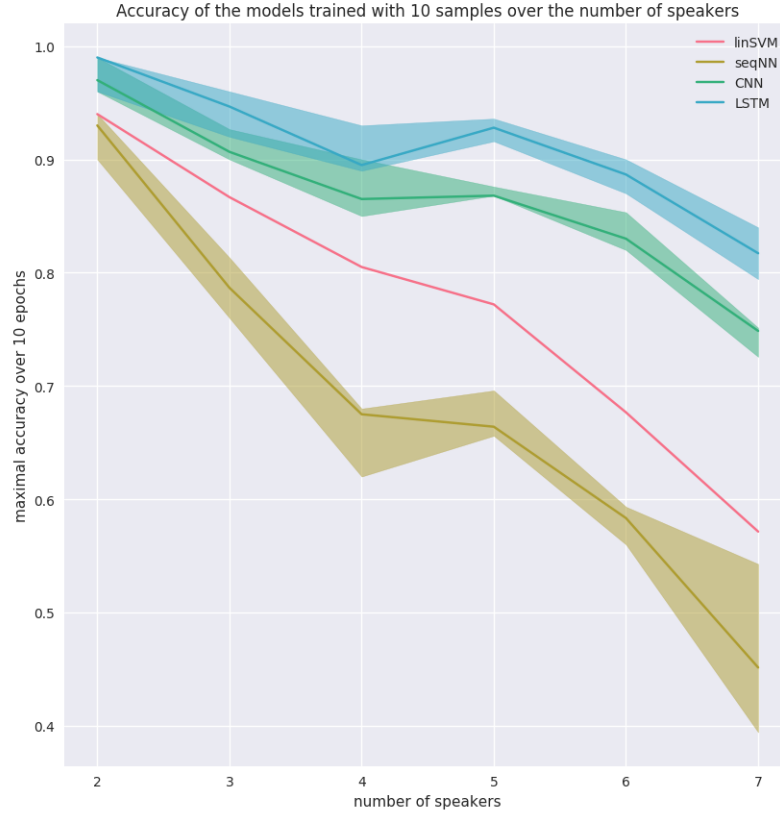


Figure 4

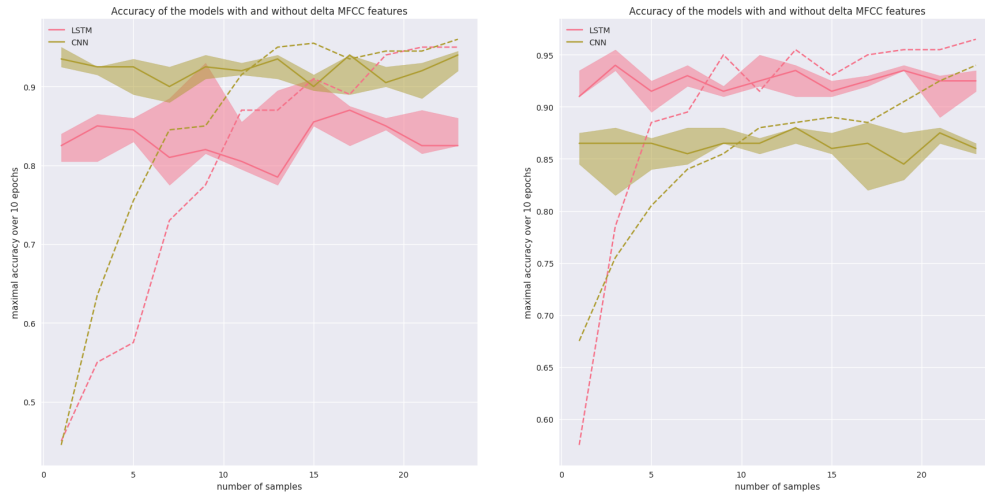
152 To put it in a nutshell, the models have all different performances. The
 153 linear Support Vector machine seemed to be the less volatile due to its de-
 154 terministic training process whereas the performances of the neural networks
 155 were more variables. The SVM could also train on less samples than the
 156 deep learning approaches which needed a more significant training set to be
 157 accurate.

158 In our final build, we implement a LSTM model as it seemed to be the best
 159 performing one and we used it as the reference method for the DSP testing.

160 2.2. DSP results

161 We implemented several DSP augmentations to our system to try to
 162 achieve a better classification accuracy. First, we tried to append the delta
 163 Mel coefficient matrix to our vector of features. We got some interesting
 164 results as seen on the graphs below. we plotted the accuracy to distinguish
 165 between four people, two male two female, of our two best models over the
 166 number of training samples. The dotted lines represent the results without
 167 the delta coefficients and the plain bold lines are the median of three same
 168 experiment with the delta inputed as features. The testing was done with 50
 169 test samples per speaker.

170



(a) Accuracy without sample extension (b) Accuracy with sample extension

Figure 5: Accuracy with vs without delta coefficient

171 We can see that the delta coefficients greatly increase the accuracy of
 172 the models for very low number of samples. CNN without delta nor sample
 173 extension is below .7 of accuracy for a set of training sample less than 5 per
 174 speaker whereas with delta features added it is above .9. LSTM without
 175 delta nor sample extension is below .8 of accuracy for less than 9 training
 176 samples per speaker whereas it is always above .8 with the delta features
 177 included. With sample extension, the results are similar but we notice once

178 more that LSTM starts performing better from the 13th sample. CNN with
179 delta also decreases in accuracy: its median score was above .9 without sam-
180 ple extension and between .8 and .9 with sample extension. It follows that
181 those added features describe some voice features different from the Mel co-
182 efficients and which characterize well the speaker. However, we can see that
183 the models trained without delta coefficients tend to be better performing as
184 the number of samples increases. Indeed, in figure (a) CNN without deltas
185 outperform CNN with delta after sample number 12 and LSTM after sample
186 10. Our theory is that our neural network architectures are not optimal and
187 that some over-fitting occurs if the vector of features is too big. We can also
188 notice one more time that without sample extension, CNN is having better
189 performances whereas LSTM is better with that feature enabled.

190
191 Another interesting result is the effect of noise reduction. Using the
192 SciPy library[6], we designed a finite impulse response filter using the Parks-
193 McClellan algorithm to filter out ambient noise. Human vocal frequency most
194 commonly lies between 300 and 3000 Hz [8], so our filter acts as a bandpass
195 that attenuates frequencies outside this range with a passband transition
196 width of 250 Hz. However, while the audio output from this filter sounds
197 noticeably better than the original to our ears, our tests show that the inte-
198 gration of such a filter dramatically reduces the accuracy of our models.

199

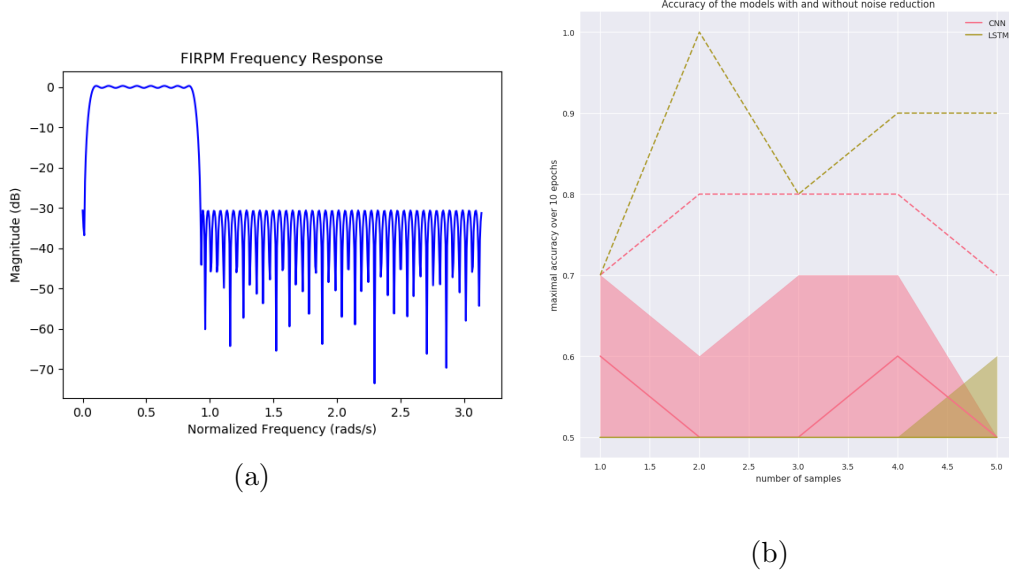


Figure 6

200 The dotted lines represent the models trained with unfiltered voice sam-
 201 ples, and the straight lines represent the models trained with noise reduced
 202 voice samples. We think that, while trying to filter out the ambient noise,
 203 the filter might have removed important features that the model uses to
 204 distinguish users. Additionally, if our test samples were taken in different
 205 conditions (hence different ambient noise), the reduction might have shown
 206 to be more efficient. We also notice that what our ears seem to like might
 207 not necessarily be aligned with what the models would require for optimal
 208 predictions.

209

210 Finally, we tested down sampling with the goal of reducing computation
 211 time, as motivated in the objectives part. A simple down sampling by a
 212 factor of 2 gave us an important decrease in accuracy, especially for small
 213 number of user voice samples.

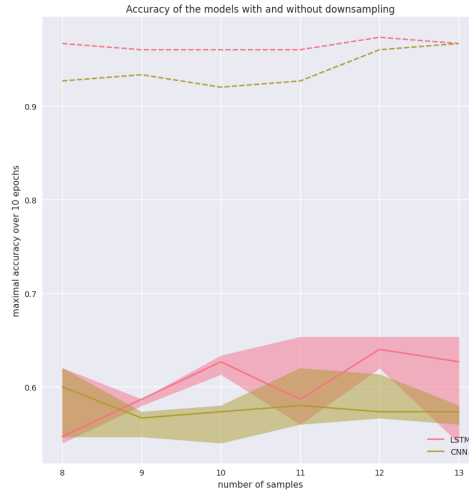


Figure 7

214 In figure 6, the dotted lines represent the accuracy of the prediction with
 215 the original voice inputs. The full lines show the accuracy with down sampled
 216 voice inputs. We can clearly see that keeping the sampling rate unconverted
 217 leads to better predictions (above 90 percent for original inputs, compared to
 218 60 percent for down sampled inputs). We assume this loss in accuracy is due
 219 to a loss in characteristics of the human voice samples when we down sample
 220 the signal. Moreover, the MFC vectors here have half as many coefficients
 221 compared to original signals. The smaller quantity of information provided
 222 to the training models might also have an impact in the overall decrease in
 223 prediction accuracy. Because of this, we have decided to discard sample rate
 224 conversion in our final system.

225 3. Conclusion

226 Our speaker recognition system demonstrated decent accuracy and ac-
227 ceptable latency when provided with a set of ten 3-second voice samples.
228 Our results were particularly good using the public dataset of voice files. We
229 discovered that LSTM was the best choice out of the four models we tested,
230 and optional features like delta features and sample extension were situation-
231 ally useful. However, additional modifications would be needed to achieve
232 similar performance in practical environments.

233
234 Complications unique to real-world environments turned out to be among
235 our greatest challenges while trying to improve the performance of our sys-
236 tem. In particular, we’ve learned that the reliability of MFC coefficients
237 appear to be highly susceptible to ambient noise, and distinguishing similar-
238 sounding voices can be difficult for a classifier. Extensions to our system
239 should focus on improving its real-world performance, either in accuracy
240 (noise reduction, classification) or usability (latency, samples required, unique
241 users supported).

242
243 There are several modifications to our system we think are worth pur-
244 suing. For example, we saw that incorporating delta features in the MFC
245 coefficients improved performance when training with few samples. Because
246 of this, we think that perhaps incorporating delta-delta features, also known
247 as acceleration features, might boost accuracy even more.

248
249 Also, since our bandpass noise reduction filter did not help the accuracy of
250 our model, we’ve learned that noise-reduced audio signals that sound better
251 than the original to our ears are not necessarily better for training classifiers
252 since important vocal features may have been too heavily attenuated. A gen-
253 tler FIRPM filter or more comprehensive tuning might resolve this, but we
254 suggest trying an adaptive filter to take advantage of the fact that real-world
255 samples will likely have the same ambient noise throughout the recording.

256
257 With modifications like these as well as careful consideration regarding
258 the trade-offs between accuracy and usability, our speaker recognition system
259 could very well have a level of robustness ready for practical applications.

260 4. Bibliography

- 261 [1] J. K. MacCallum, A. E. Olszewski, Y. Zhang, J. J. Jiang, Effects of
262 low-pass filtering on acoustic analysis of voice, *J Voice* 25 (2011) 15–20.
263 20346621[pmid].
- 264 [2] R. B. Miller, Response time in man-computer conversational transac-
265 tions, in: *Proceedings of the December 9-11, 1968, Fall Joint Computer*
266 *Conference, Part I, AFIPS '68 (Fall, part I)*, ACM, New York, NY, USA,
267 1968, pp. 267–277.
- 268 [3] S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural Com-*
269 *put.* 9 (1997) 1735–1780.
- 270 [4] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion,
271 O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vander-
272 plas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay,
273 Scikit-learn: Machine learning in Python, *Journal of Machine Learning*
274 *Research* 12 (2011) 2825–2830.
- 275 [5] F. Chollet, et al., Keras, 2015.
- 276 [6] A. Torfi, SpeechPy: Speech recognition and feature extraction, 2017.
- 277 [7] J. Kominek, A. W. Black, V. Ver, CMU Arctic Databases for Speech
278 Synthesis, Technical Report, 2003.
- 279 [8] Y. Hu, P. C. Loizou, Subjective comparison of speech enhancement al-
280 gorithms, in: *2006 IEEE International Conference on Acoustics Speech*
281 *and Signal Processing Proceedings*, volume 1, pp. I–I.