



A New Front-End for Classification of Non-Speech Sounds: A Study on Human Whistle

Mahesh Kumar Nandwana, Hynek Bořil, John H. L. Hansen*

Center for Robust Speech Systems (CRSS)
Erik Jonsson School of Engineering and Computer Science
University of Texas at Dallas, Richardson, Texas, USA
{mahesh.nandwana, hynek, john.hansen}@utdallas.edu

Abstract

Speech/non-speech sound classification is an important problem in audio diarization, audio document retrieval and advanced human interfaces. The focus of this study is on the development of spectral and temporal acoustic features for speech/non-speech sound classification based on production differences in speech versus whistle. Seven time- and frequency-domain based features are investigated. Performance of the proposed feature set for the task of speech/whistle classification is evaluated at frame level. This evaluation utilizes support vector machine (SVM) models and Gaussian mixture models (GMM) for back-end classifiers. At the frame-level, the proposed front-end fusion gives an absolute performance gain of +15.0% and +3.1% over MFCC with SVM and GMM based classifiers, respectively. This research will benefit the development of intelligent speech interfaces for identification, recognition, and speech coding, as a pre-processing step for real world audio streams.

1. Introduction

Human sounds produced via the oral cavity can be classified into two broad categories: i) speech and ii) non-speech. Non-speech sounds include vocalizations such as: scream, whistle, cough, laugh, snore, sneeze, hiccups, etc. Effective classification of non-speech sounds is a necessary preprocessing step for robust speech and speaker recognition, audio indexing and diarization, as well as other applications. This study focuses exclusively on classification of human whistles which are a part of the general class of non-speech sounds.

Performance of speech systems for automatic speech recognition (ASR) and speaker and language identification degrades significantly as soon as they encounter a mismatch between train and test conditions. Such a mismatch can be introduced either by speaker dependent factors or environment/hardware dependent factors. Speaker dependent factors include non-speech sounds [1], vocal effort and speaking styles [2, 3, 4], speech under stress [5], Lombard effect [6, 7], whereas environment/hardware dependent factors include microphone characteristics, room acoustics, channel mismatch, etc.

Human whistle is produced by controlling the stream of air flow generated via lungs and passing through the oral cavity. Here, the oral cavity works as a resonant chamber. Whistling can be used to get attention, to call a human or a pet, or to carry a melody. In general, human whistle can be classified into many categories such as pucker whistles, finger whistles, teeth whistles, bird whistles, warble whistles, roof whistles, etc. However, for this work, we only consider pucker whistles.

*This project was funded by AFRL under contract FA8750-12-1-0188 and partially by the University of Texas at Dallas from the Distinguished University Chair in Telecommunications Engineering held by J.H.L. Hansen.

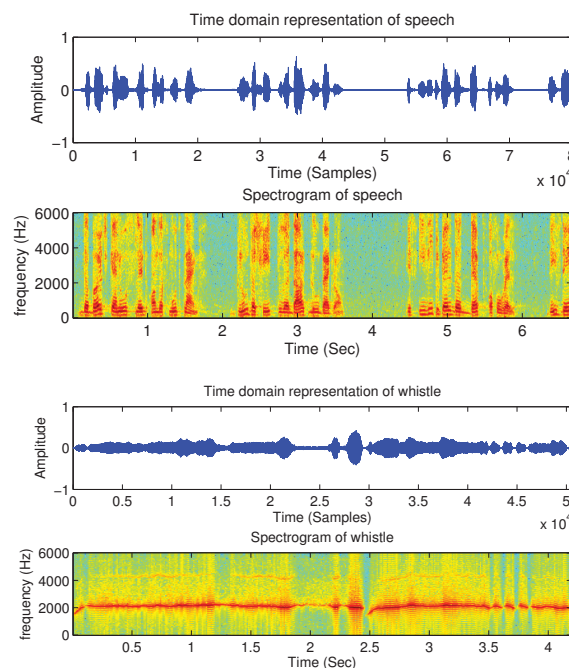


Figure 1: Time and frequency domain representation of speech (top) and human whistle (bottom).

Pucker whistle is the most common form of human whistling. Pucker whistles are produced by curving the tongue inside the oral cavity such that the top of the tongue touches the roof of the oral cavity, where the tip of the tongue should be downwards to create turbulence followed by blowing out, or sucking air from mouth. Different resonances can be produced by changing the shape of the tongue and position of the jaw.

Human whistles are signal tone sounds that contain a single dominant frequency. From the spectrogram shown in Fig. 1, we can clearly see there is a harmonic line representing the presence of a dominating frequency in the spectrogram of human whistle, whereas speech has simultaneous occurrence of multiple frequencies. Automatic speech recognizers, speaker verification systems, diarization systems, and other speech-oriented engines are prevalently trained on clean speech and when exposed to whistle sounds, their performance tends to deteriorate.

Speech/non-speech classification has been an active area in the domain of auditory scene analysis [8]. Research have considered non-speech sounds such as scream [9, 10], snore [11], and various environmental sounds [12]. However, the only work related to human whistle processing is by Nilsson [13], where the study considered 20 subjects and included frequency analysis in

noisy environments. The lack of research in the area of human whistle processing suggests that further research is needed. The applications of whistle classification can be found in the following areas:

- *Audio information retrieval*: searching whistles in movies and radio/TV shows;
- *Audio indexing*: indexing whistles in audio streams;
- *Music technology*: automatic melody extraction;
- *Environment classification*: steam whistles and train whistles classification; environmental sniffing;
- *Speech technology*: impact on speech and speaker recognition systems.

In this study, we first analyze and compare spectral and temporal properties of speech and whistle segments. Based on the analysis, we propose a set of features for speech/whistle classification. The features are evaluated side-by-side with traditional mel frequency cepstral coefficients (MFCC) using Gaussian mixture model (GMM) and support vector machine (SVM) based classifiers. Finally, we demonstrate the benefits of the proposed classification scheme on an example of speaker verification.

2. Corpus

At present, there is no publicly available corpus for human-whistle research, especially none that would allow for studies of whistle in the context of speech or speaker recognition. The corpus used in this study was collected at the University of Texas at Dallas. The recordings were captured at a 44.1 kHz sampling rate in an ASHA certified single walled sound booth using a table top microphone.

A total of 30 subjects (17 males and 13 females) participated in the corpus collection. Subjects were native as well as non-native English speakers. Each subject participated in a single recording session divided into three parts. In the first part, the subjects read the first 11 lists (110 phonetically balanced sentences) from the IEEE recommended set of phonetically balanced sentences [14]. In the second part, the subjects produced spontaneous speech while answering 6 questions. In the third part, they were asked to whistle in a sequence of audio captures. The subjects were asked to capture maximum variability in their whistling style and were not instructed to imitate any particular song or melody. The corpus details are given in Table 1. Here, each sample is defined as a 10 seconds long speech or whistle audio chunk. The samples incorporate silence segments that naturally occurred during the speech/whistle production.

3. Front-End for Speech/Whistle Classification

In this section, speech and whistle samples are analyzed in terms of their spectral and temporal properties. Based on the observations, a set of features for speech/whistle classification is proposed.

3.1. Analysis of Speech and Whistle

In voiced portions of speech, an air flow pushed from the lungs causes vibration of the vocal folds. This vibration serves as excitation for the vocal tract resonances. In unvoiced portions of speech (unvoiced consonants or whispered speech) the glottis is

Class	# Samples	Total Duration (min)
Speech	1933	322
Whistle	247	45.6

Table 1: Whistle corpus details.

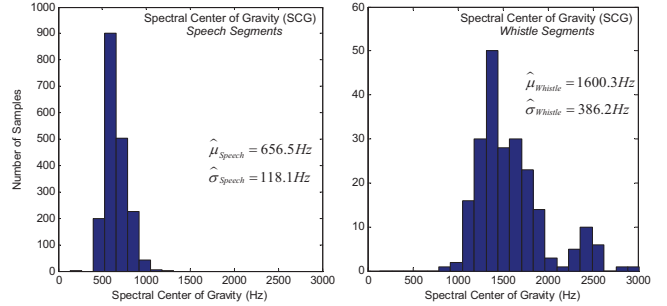


Figure 2: Distribution of spectral center of gravity (SCG) in speech and whistle samples drawn from the Whistle corpus.

kept open and a turbulent flow produced by the passing air serves as a source for the articulators. As a result, voiced portions of speech can be represented in the spectral domain as a series of glottal waveform harmonics weighted by the transfer function of the vocal tract, and unvoiced portions as a weighted spectrum of noise. While exhibiting different spectral slopes (unvoiced speech having a flatter tilt), the spectrum of both voiced and unvoiced speech spans the whole frequency range considered in traditional speech processing schemes. In comparison, the production of whistle results in a sharp peak in the spectrum representing the fundamental frequency, with a potential presence of weak higher harmonics and a noise-like spectral component representing the unvoiced sound created by the friction of the air stream passing through the constrictions created by the tongue, lips (and possibly fingers) while whistling (see Fig. 1). For this reason, it can be expected that the distribution of spectral energy in speech and whistle will display significant differences. For the spectral analysis purposes, as well as with respect to the objective of automatic segment classification later in the text, we choose parametric representations of the amplitude spectra.

Spectral center of gravity (SCG), representing the ‘center of mass’ of the power spectrum, and spectral energy spread (SES), which represents the standard deviation of the spectral energy distribution from SCG [15], were extracted from speech and whistle samples of the Whistle corpus and analyzed (see Fig. 2 and 3). It can be seen that the SCG distribution of speech is sharper and centered at a significantly lower frequency compared to whistle. The whistle SCG distribution directly reflects the range of whistle fundamental frequencies produced by the subjects. Due to the variety in the choice of a whistled melody and pitch, the overall SCG distribution extracted across all subjects displays larger variation compared to speech. As expected, SES distributions in Fig. 3 show the opposite trend compared to SCG – the energy spread around SCG within individual sam-

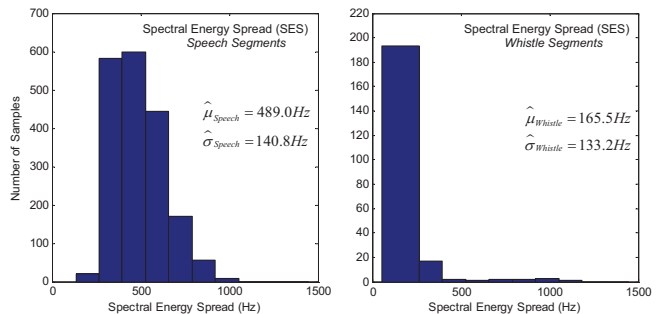


Figure 3: Distribution of spectral energy spread (SES) in speech and whistle samples drawn from the Whistle corpus.

ples is much wider for the ‘broad’ speech spectrum compared to sharp whistle spectra. The considerable differences in SCG and SES distributions of speech and whistle suggest these two parameters could be valuable for the speech/whistle classification.

In the next step, we analyze pitch contours in speech and whistle parameterized by so called pitch patterns. While the pitch of modal speech and whistle clearly occupy different frequency regions, as already reflected in the above mentioned spectral parameters, it may be interesting to see whether there are also differences in their time trajectories independent of their absolute positioning in frequency. For this purpose, we perform a pitch pattern analysis introduced in [16]. In this method, after traditional frame-level pitch extraction, continuous voiced sections are median-filtered and processed by a regression analysis using a sliding window. Voiced sections shorter than the window are dropped. A straight line is fit into the F_0 contour within the window by means of linear regression. If the regression line is steep enough to cross the whole ‘threshold’ frequency band within the length of the analysis window, the segment is assigned a rising/falling pattern element; otherwise a flat pattern is assigned. In the next step, N gram statistics of the consecutive pitch pattern elements representing the pitch contour are computed.

In our study, the RAPT cross-correlation algorithm [17] implemented by WaveSurfer [18] is used for pitch extraction, and a window length of 50 ms with a 100 % skip rate, together with a 5 Hz threshold frequency band are used for the pitch pattern extraction. First, frequencies of unigram patterns were extracted for speech (*up* – 22.3 %; *flat* – 45.2 %; *down* – 32.4 %) and whistle (*up* – 40.2 %; *flat* – 28.1 %; *down* – 31.7 %). It can be seen that while the *down* pattern occurred with similar frequency in speech and whistle, the *flat* pattern was chosen by the subjects more frequently in speech and *up* in whistle. Subsequently, frequencies of pitch pattern bigrams were analyzed (see Fig. 4). It can be seen that similarly as for unigrams, a *flat-flat* bigram dominates speech pitch contours and *up-up* is dominating in whistle. The overall distribution of pitch pattern bigrams is more uniform in whistle than in speech, suggesting higher variability of the whistled pitch contours. This suggests the subjects complied with the instructions to produce a variety of whistle tones rather than just stationary whistles.

The pitch pattern models discussed here show a good potential to contribute to speech/whistle classification as additional features, especially considering the fact that prosody-related content is known to be relatively consistent within a spoken language (pitch patterns were successfully leveraged for example in language identification [19]). On the other hand, pitch pattern histograms are likely to be consistent within individual types of whistle (e.g., single tone versus melodic whistle) while varying across whistle types. To avoid a bias in our study on automatic speech/whistle classification through having an *a priori* knowledge of the whistle type (melodic), we refrain from utilizing pitch patterns in the following experiments.

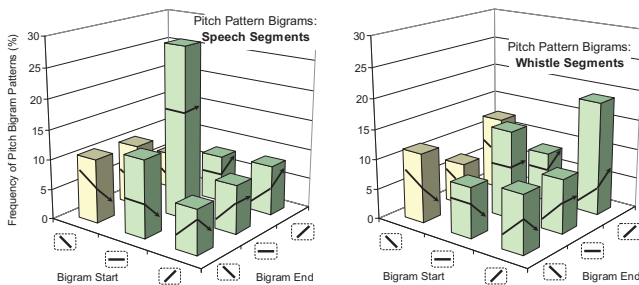


Figure 4: Frequency of bigram pitch patterns in speech and whistle.

3.2. Features for Speech/Whistle Classification

Based on the fundamental differences in speech and whistle spectra and pitch observed in the previous section, we propose to use a combination of time-domain and spectral-based parameters as features for speech/whistle classification. The features include: zero-crossing rate (ZCR), spectral center of gravity (further denoted as spectral centroid – SC), spectral energy spread (further denoted as ‘SS’), spectral crest factor (SCF), spectral decrease (SD), spectral kurtosis (SK), and spectral skewness (SSk). Several of these features have previously been considered for general audio content analysis [20], speaking style classification [21], and in emotion and cognitive load classification [22]. Since most of these parameters are widely used in the community, we will provide only a short summary of their properties. *Zero Crossing* represents the rate of signal sign changes within a segment and is popular for speech activity detection. *Spectral Crest Factor* is represented by the ratio of the maximum of the magnitude spectrum over the sum of the magnitude spectrum. It is a measure of the tonalness of an audio signal. *Spectral Decrease* is the measure of the steepness of the decrease of the spectral envelope over frequency. *Spectral Kurtosis* is a measure of the spectral ‘peakedness’. *Spectral Skewness* measures the symmetry of the distribution of the spectral magnitude values around their arithmetic mean.

To complement the analysis from Sec. 3.1, means and variances of all parameter distributions are summarized in Table 2. For classification purposes, the time and frequency domain features are concatenated into a seven-dimensional feature vector where each dimension (feature) is mean and variance normalized on an sample level.

4. Experiments and Results

In this section, effectiveness of the proposed front-end for speech/whistle classification is evaluated. Performance evaluation is performed at the frame-level using SVM and GMM based classifiers. About 6K frames are used for training and 31500 frames for testing. The training is speaker and gender independent and a balanced amount of whistle and speech frames are provided in both the training and the test set. The classification features are extracted using a Hamming window of 25 ms with a skip rate of 10 ms. Classification results are reported in terms of accuracy, which is the ratio of correctly classified labels over the total number of labels.

In the frame-level classification, the goal is to label whether a frame belongs to speech or whistle. When utilizing individual features (one feature at a time), classification accuracies ranged between 72 % and 78 %. We also considered various combinations of the individual features (2 to 6) from the proposed feature vector, but with a few exceptions, the performance decreased. Therefore, we conclude that all seven features contribute to the classification task. As can be seen in Table 3, with one exception, the classifier utilizing the proposed feature vector outperforms the MFCC baseline both for the SVM and GMM setups. Me-

Feature	Speech		Whistle	
	mean	std	mean	std
SC	656.5	118.1	1600.3	386.2
SCF	0.11	0.06	0.19	0.07
SD	0.05	0.04	0.01	0.05
SK	329.29	241.47	580.54	226.03
SSk	3.77	1.46	5.54	1.32
SS	489.0	140.8	165.5	133.2
ZCR	0.21	0.14	0.36	0.10

Table 2: Distributions of selected acoustic features

Front-End	#Dim	No Filter		3pt Median Filter		5pt Median Filter		7pt Median Filter	
		SVM	GMM	SVM	GMM	SVM	GMM	SVM	GMM
MFCC	12	75.4	90.0	75.2	91.6	75.2	94.2	75.1	94.6
Proposed	7	84.4	91.4	87.1	93.3	88.8	94.6	89.9	94.2
Fusion	19	84.7	94.8	87.4	97.0	89.2	96.9	90.1	97.7

Table 3: Frame-level classification results for different front end.

dian filtering of the frame-level decisions further improves the classification performance. Finally, fusion of traditional MFCC features with the proposed feature vector is evaluated and as can be seen in the last row of the table, further improves the classification accuracy for all setups. It can be observed that in general, the GMM-based classification outperforms the SVM one.

5. Impact on Speaker Verification System

5.1. GMM-UBM Framework

Our earlier study has considered human scream, which represents another class of non-speech sounds, and reported its impact on the speaker recognition systems [1]. In this section, we investigate the impact of human whistles on the performance of a speaker verification system. To compensate for the train/test mismatch, where the speaker models are trained on speech segments and the test tokens are a mixture of speech and whistle, we utilize the speech/whistle classifier proposed in the previous section. The classifier is applied to detect and drop frames containing whistle so the test frames scored against the speaker models would contain pure speech. From the corpus of 30 speakers, we use 25 speakers for the speaker verification experiments and the remaining 5 for training a speech/whistle classifier. The entire scheme is depicted in Fig. 5.

We use a GMM-UBM framework for speaker recognition. A universal background model (UBM) is constructed using ‘‘CRSS-4English-14 corpus’’ which has a total of 454 speakers from four different dialects of English. This corpus was collected in identical conditions (microphone channel, recording room) as the whistle corpus. For UBM training, a subset of 225 speakers are selected with a total of 1801 sentences. A speaker specific maximum a posteriori (MAP) adapted Gaussian mixture model (GMM) is obtained from the UBM for each of the trained speakers [23]. The test files are scored against the UBM and speaker specific GMM and resulting scores were used to obtain overall system accuracy. Performance is evaluated by computing equal error rate (EER) for the ensemble of trials.

5.2. Front-End Processing

To observe the effect of whistle on speaker verification, audio files consisting of sequences of alternated speech and whistle segments were constructed. The ratio of speech to whistle seg-

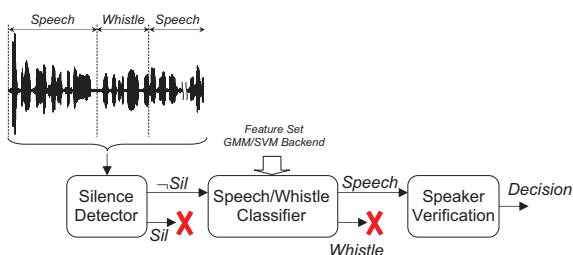


Figure 5: Proposed compensation scheme for speech/whistle mismatch in speaker verification.

Test Condition	Baseline	With Whistle Detection			
		Proposed		Proposed+MFCC	
		SVM	GMM	SVM	GMM
Speech	10.8	---			
Speech +Whistle	21.6	16.3	18.0	15.4	15.1

Table 4: Speaker verification compensation results; EER (%).

ments within speech-whistle audio streams was seven to four.

The front-end processing includes down-sampling the data to 8 kHz. Silence frames are dropped using an energy and zero crossing rate based voice activity detector (VAD). 36-dimensional MFCC vectors containing c_1 - c_{12} static coefficients and corresponding delta and delta-delta coefficients are extracted using a Hamming window of 25 ms with a skip rate of 10 ms. Cepstral mean and variance normalization is also applied across all features at the sample (token) level.

5.3. Speaker Verification Experiment

A total of five 10-seconds long speech samples per speaker are used to produce 64 mixture GMM speaker-specific models via MAP adaptation of the UBM. A total of 62500 random trials were generated for speaker verification. The ratio of target to impostor trials was 1:24 (2500 targets and 60000 impostors). Similar to training, each test sample is also 10 seconds long. Speaker verification results are summarized in Table 4.

EER in the case of testing with clean speech samples is 10.80%. Although the absolute EER for matched training/testing speech condition can be brought further down by increasing the number of Gaussian mixture components, to limit the computational demands, 64 mixture components are used throughout all experiments here. In the case speech-whistle mixture samples, WER increases to 21.60%. It is evident that the speaker verification system cannot sustain its performance when exposed to audio containing whistle islands. Thus, to compensate for the effects of whistle, we incorporate the proposed speech/whistle classification scheme to identify and drop whistle frames from the test audio stream. From Table 4, it is clear that automatic dropping of whistle frames reduces EER by a large margin. We note that the absolute duration of the speech frames in the mixed speech and whistle test samples is about 7 seconds and hence, even after successfully removing whistle segments, the verification system has access to less speech frames than in the baseline matched ‘speech only’ task. For this reason, the WERs here are higher than the baseline ones.

6. Conclusion and Future Work

In this study, it was observed that the presence of human whistle in audio streams impacts performance of speech systems. To make the systems more robust, new features to classify speech and whistle frames were proposed. Overall gains in the speech/whistle classification were reached using the proposed features combined with two alternative classifier backends. The proposed front-end fusion provided an absolute performance gain of +15.0% and +3.1% over MFCC with SVM and GMM classifiers, respectively. A compensation scheme for non-speech mismatch reduction was also evaluated and shown to improve performance of a speaker verification system. A wider range of environmental sounds could be explored as the next step for multi-class non-speech sound classification.

7. References

- [1] M. K. Nandwana and J. Hansen, "Analysis and identification of human scream: Implications for speaker recognition," *INTER-SPEECH 2014*, pp. 2253–2257, 2014.
- [2] C. Zhang and J. Hansen, "Analysis and classification of speech mode: whispered through shouted," *INTERSPEECH 2007*, pp. 2289–2292, 2007.
- [3] E. Shriberg, M. Graciarena, H. Bratt, A. Kathol, S. S. Kajarekar, H. Jameel, C. Richey, and F. Goodman, "Effects of vocal effort and speaking style on text-independent speaker verification," *INTER-SPEECH 2008*, pp. 609–612, 2008.
- [4] S. Ghaffarzadegan, H. Bořil, and J. H. L. Hansen, "Model and feature based compensation for whispered speech recognition," in *Interspeech 2014*, Singapore, Sept 2014, pp. 2420–2424.
- [5] J. Hansen, "Analysis and compensation of speech under stress and noise for environmental robustness in speech recognition," *Speech communication*, vol. 20, no. 1, pp. 151–173, 1996.
- [6] J. Hansen and V. Varadarajan, "Analysis and compensation of Lombard speech across noise type and levels with application to in-set/out-of-set speaker recognition," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 2, pp. 366–378, Feb 2009.
- [7] H. Bořil and J. H. L. Hansen, "Unsupervised equalization of Lombard effect for speech recognition in noisy adverse environments," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 6, pp. 1379–1393, August 2010.
- [8] L. Lu, H. J. Zhang, and H. Jiang, "Content analysis for audio classification and segmentation," *IEEE Transactions on Speech and Audio Processing*, vol. 10, no. 7, pp. 504–516, Oct 2002.
- [9] M. K. Nandwana, A. Ziaei, and J. H. L. Hansen, "Robust unsupervised detection of human screams in noisy acoustic environments," in *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on*, 2015.
- [10] G. Valenzise, L. Gerosa, M. Tagliasacchi, E. Antonacci, and A. Sarti, "Scream and gunshot detection and localization for audio-surveillance systems," in *Advanced Video and Signal Based Surveillance, 2007. AVSS 2007. IEEE Conference on*, Sept 2007, pp. 21–26.
- [11] W. H. Liao and Y. K. Lin, "Classification of non-speech human sounds: Feature selection and snoring sound analysis," in *Systems, Man and Cybernetics, 2009. SMC 2009. IEEE International Conference on*, Oct 2009, pp. 2695–2700.
- [12] J. Portelo, M. Bugalho, I. Trancoso, J. Neto, A. Abad, and A. Serralheiro, "Non-speech audio event detection," in *Acoustics, Speech and Signal Processing, 2009. ICASSP 2009. IEEE International Conference on*. IEEE, 2009, pp. 1973–1976.
- [13] M. Nilsson, J. Bartunek, J. Nordberg, and I. Claesson, "Human whistle detection and frequency estimation," in *Image and Signal Processing, 2008. CISP '08. Congress on*, vol. 5, May 2008, pp. 737–741.
- [14] "IEEE recommended practice for speech quality measurements," *IEEE Transactions on Audio and Electroacoustics*, vol. 17, no. 3, pp. 225–246, Sep 1969.
- [15] H. Bořil, O. Sadjadi, and J. H. L. Hansen, "UTDrive: Emotion and cognitive load classification for in-vehicle scenarios," in *The 5th Biennial Workshop on Digital Signal Processing for In-Vehicle Systems*, Kiel, Germany, September 2011.
- [16] H. Bořil, Q. Zhang, P. Angkititrakul, J. H. L. Hansen, D. Xu, J. Gilkerson, and J. A. Richards, "A preliminary study of child vocalization on a parallel corpus of US and Shanghai toddlers," in *INTERSPEECH 2013*, Lyon, France, 2013, pp. 2405–2409.
- [17] D. Talkin, *Speech Coding and Synthesis*. Amsterdam, Netherlands: Elsevier, 1995, ch. A Robust Algorithm for Pitch Tracking (RAPT). W.B. Kleijn and K.K. Paliwal (Eds.), pp. 495–518.
- [18] K. Sjolander and J. Beskow, "WaveSurfer – an open source speech tool," in *Proc. of ICSLP'00*, vol. 4, Beijing, China, 2000, pp. 464–467.
- [19] H. Bořil, Q. Zhang, A. Ziaei, J. H. L. Hansen, D. Xu, J. Gilkerson, J. A. Richards, Y. Zhang, X. Xu, H. Mao, L. Xiao, and F. Jiang, "Automatic assessment of language background in toddlers through phonotactic and pitch pattern modeling of short vocalizations," in *Workshop on Child Computer Interaction (WOCCI)*, Singapore, September 2014.
- [20] A. Lerch, *An introduction to audio content analysis: Applications in signal processing and music informatics*. John Wiley & Sons, 2012.
- [21] H. Bořil, P. Fousek, and H. Höge, "Two-stage system for robust neutral/Lombard speech recognition," in *Proc. of Interspeech'07*, Antwerp, Belgium, 2007, pp. 1074–1077.
- [22] H. Bořil, O. Sadjadi, T. Kleinschmidt, and J. H. L. Hansen, "Analysis and detection of cognitive load and frustration in drivers speech," in *Interspeech'10*, Makuhari, Chiba, Japan, September 2010, pp. 502–505.
- [23] D. A. Reynolds, T. F. Quatieri, and R. B. Dunn, "Speaker verification using adapted Gaussian mixture models," *Digital Signal Processing*, vol. 10, no. 1–3, pp. 19–41, 2000.