

Project no: FP6 034362

ACORNS
Acquisition of Communication and Recognition Skills

Instrument: STREP

Thematic Priority: Information Society Technologies

D1.2: Modules for a) augmentation of standard spectral features with a stream of milli-second and decisecond features and evaluation on specific phone classification tasks and b) feature selected by sensitivity-analysis method

Due date of deliverable: 2008-12-01

Actual submission date: 2008-12-01 (latest version)

Start date of project: 2006-07-01

Duration: 36 Months

Organisation name of lead contractor for this deliverable: KTH

Revision: 1.1

Project co-funded by the European Commission within the Sixth Framework Programme 2002-2006		
Dissemination Level		
PU	Public	X
PP	Restricted to other programme participants (including the Commission Services)	
RE	Restricted to a group specified by the consortium (including the Commission Services)	
CO	Confidential, only for members of the consortium (including the Commission Services)	

Table of Contents

1	OVERVIEW	6
2	AUTOMATIC VOICE ONSET TIME ESTIMATION FROM REASSIGNMENT SPECTRA	8
2.1	INTRODUCTION	8
2.2	SPECTRAL REASSIGNMENT	10
2.3	PROPERTIES OF THE VOICE ONSET TIME	11
2.4	DATA SETS	12
2.4.1	THE "FORCED" DATA SET	12
2.4.2	THE "FREE" DATA SET	13
2.4.3	THE "MANUAL" DATA SET	13
2.4.4	THE "TEST" SET	13
2.5	THE VOT ESTIMATION ALGORITHM	14
2.5.1	DETECTION OF PLOSIVE SEGMENTS	14
2.5.2	BURST ONSET DETECTION	15
2.5.3	START OF PERIODICITY	15
2.5.4	DISCUSSION	16
2.6	EXPERIMENTS	16
2.6.1	ALGORITHM PERFORMANCE FOR PHONETIC STUDIES	16
2.6.2	ALGORITHM PERFORMANCE FOR AUTOMATIC SPEECH RECOGNITION	18
2.6.3	ESTIMATED VOTs	19
2.6.4	VOT AS A FEATURE FOR AUTOMATIC SPEECH RECOGNITION	21
2.7	CONCLUSIONS	24
3	USING PROSODY WITH NMF	27
3.1	INTRODUCTION	27
3.2	RHYTHM VECTORISATION	27
3.3	RHYTHM VQ RESULTS	29
3.4	PITCH VECTORISATION	30
3.5	PITCH VQ RESULTS	32
3.6	CONCLUSIONS	32
4	FEATURE SELECTION BASED ON KNOWLEDGE OF THE AUDITORY SYSTEM	33
4.1	INTRODUCTION	33
4.2	PRELIMINARIES	34
4.2.1	DISTORTION ANALYSIS	34
4.2.2	VAN DE PAR AUDITORY MODEL	35
4.3	A METHOD FOR RATING THE PERCEPTUAL SIGNIFICANCE OF FEATURES	35
4.3.1	FEATURE AND SPEECH: A LINEARIZED RELATION	35
4.3.2	THE SENSITIVITY MATRIX IN THE FEATURE DOMAIN	36
4.4	APPLICATION TO SPEECH RECOGNITION	37
4.4.1	THE A MATRIX FOR MFCCs	37
4.4.2	THE SELECTION PROCESS	38

4.5	EVALUATION	38
4.5.1	RANGE OF LINEARIZATION	38
4.5.2	SPEECH RECOGNITION EXPERIMENTS	39
4.6	CONCLUSIONS	39

Abstract

This report describes the progress within the ACORNS project towards augmentation of the standard spectral features used for speech recognition with millisecond and decisecond features and towards feature selection based on knowledge of the human auditory system. Our work on millisecond features has resulted in an algorithm that automatically estimates the voicing onset of plosives. Our work towards defining decisecond features consists of a study on the efficacy of including measures of prosody (rhythm and pitch movement). We find that the measures are useful at the initial learning stage. Our work on feature selection led to an algorithm that selects features based on the ability of the features to describe the components of speech that are most clearly perceived. Experimental results confirm effectiveness of this generic strategy.

Chapter 1

Overview

This report is the description component of deliverable D1.2 of the ACORNS project. As listed in the Annex, the deliverable was aimed to consist of “Modules for a) augmentation of standard spectral features with a stream of milli-second and decisecond features and evaluation on specific phone classification tasks and b) feature selection by sensitivity-analysis method (software and report)”. Part a) of the deliverable is associated with Task 2 (“phone-class specific features) and part b) with Task 1 (“distortion-based approach”).

With respect to Task 2, our work towards defining milli-second features has taken the form of an algorithm that automatically estimates the voicing onset of plosives. The work towards defining and decisecond features is a study on the efficacy of including measures of prosody (rhythm and pitch movement) in the feature set used in the ACORNS project. With respect to Task 1, we developed an algorithm to select from a larger set of features a subset of features based on the ability of the subset to describe the audible components of the signal.

This report part of the deliverable consists of three chapters, each describing one of the fore-mentioned topics. Chapters 2 and 4 are based on papers that are currently under review.

Chapter 2 describes the algorithm to automatically estimate the voice onset time (VOT) of plosives. The VOT is the time delay between the burst onset and the start of periodicity when it is followed by a voiced sound. Since the VOT is affected by factors like place of articulation and voicing it can be used for inference of these factors. The algorithm uses the reassignment spectrum of the speech signal, a high resolution time-frequency representation which simplifies the detection of the acoustic events in a plosive. The performance of our algorithm is evaluated on a subset of the TIMIT database by comparison with manual VOT measurements. On average, the difference is smaller than 10 ms for 76.1% and smaller than 20 ms for 91.4% of the plosive segments. We also provide analysis statistics of the VOT of /b/, /d/, /g/, /p/, /t/ and /k/ and experimentally verify some sources of variability. To illustrate its use, we integrate the automatic VOT estimates as an additional feature in an HMM-based speech recognition system and show a small but statistically significant improvement in phone recognition rate. The new features are ready for integration in the ACORNS system.

Chapter 3 describes our study towards defining decisecond features based on prosody. The specific aim of the study was to investigate whether the addition of prosodic cues, rhythm and pitch, would increase word detection accuracy for the nonnegative-matrix factorization (NMF) approach. It was hypothesised that the addition of prosodic cues to the input stream would increase the accuracy at a faster rate. As hypothesised, the addition of prosodic cues as an aid for word detection helped raise accuracy results during the early learning period. Rhythm had more of an impact than the use of the pitch contour with accuracy almost 2% better than the baseline during the first 100 utterances. After

this period the accuracy tends towards the baseline for both cues. Calculating the pitch contour with dynamic programming smoothing slightly enhanced the results, but at the expense of computational complexity. Since the method aids only in the initial learning phase, we probably will not use these features in the ACORNS system.

Chapter 4 describes our algorithm to find a good subset of features for recognition from a larger set, using only knowledge of the human auditory system as a measure. The underlying assumption of our work is that the human auditory system is effective at extracting relevant information from the speech signal. We use a psycho-acoustic model to perform a sensitivity analysis on speech, based on a distortion measure. The method eliminates the dependency of the feature set to the speech recognition system used, and results in a generic set of good features. We evaluated the selected feature subsets on a real speech recognizer. The results confirm that knowledge of the human auditory system forms a good basis for selecting a subset of features from a larger set for the purpose of speech recognition.

While the feature selection method of Chapter 4 can be used within ACORNS to limit the number of features, this is not a major concern for project. Rather, the method should be seen as a first step towards a method that uses auditory-knowledge to improve existing features and define new features. Work towards this goal is currently in progress and will be described in deliverable D1.3. The new features will fit naturally within the ACORNS concept: the features under development reflect the innate perception of a baby (or human in general), and are not based on the recognition performance of some of some automatic speech recognition system.

The software for the three components forms part of the software that is available to the partners of the ACORNS project. As indicated in the project Annex, the software is not publicly available.

Chapter 2

Automatic Voice Onset Time Estimation from Reassignment Spectra

V. Stouten, H. Van hamme (KU Leuven)

2.1 Introduction

State-of-the-art automatic speech recognition (ASR) systems typically use a sliding window with a length of about 30 ms and a shift of about 10 ms to extract features such as Mel Frequency Cepstral Coefficients (MFCCs) from the acoustic waveform of the speech signal. However, plosives also exhibit distinctive acoustic events at a finer time scale. Typically, the closure interval ends in an abrupt increase in acoustic energy across the frequency range. The release interval is measured from this burst onset to the start of periodicity or to the onset of noise or silence. The duration of the release interval is then called voice onset time or VOT in case periodicity is present. These events can be as short as a few milliseconds. Nevertheless, they contain potentially important information on the plosive identity which is lost when a sliding window of the mentioned size is used. The subsampling caused by the 10 ms frame shift is too slow to accurately represent the timing of the events that define the release interval and the window length is too large to accurately resolve the very distinct phases of the plosive. The length of the sliding window and the frame rate that are used by today's ASR systems are a global compromise on all phones, involving e.g. effects of the variance of the spectral estimator, the trade-off between temporal and frequency resolution as dictated by the Heisenberg inequality, the data rate and the modelling constraints imposed by the subsequent acoustic modelling techniques such as Hidden Markov Models (HMMs).

Recently, there has been considerable interest in supplementing ASR systems with information that is lost during frame-based front-end processing or that is difficult to model with popular methods such as HMMs or (hybrid) Multilayer Perceptrons [1]. For instance, the phone or state duration distributions implied in an HMM match poorly with actual distributions measured on speech. In general, timing at different scales is poorly modeled in traditional ASR systems. Minor ASR accuracy improvements were found with phone duration models by [2], but the elapsed time between acoustic events at the smallest scale such as in the current VOT study, or at larger scales such as for prosodic breaks seem to be difficult to integrate in an ASR system. The work reported in [1] also illustrates that the exploitation of speech attributes like the VOT is a substantial piece of research.

The emphasis of this paper is on the automatic measurement of the VOT itself including an ac-

curacy analysis. The fact that VOT is not a frame-synchronous feature but that it is measured at the phone level and that it is only relevant for a subset of phones makes direct integration in an HMM architecture difficult. However, though we realize that this is a suboptimal approach, we will illustrate the usefulness of the VOT feature by rescored phone lattices generated by an HMM-based phone recogniser. Newer statistical modelling frameworks such as graphical models [3] probably offer additional opportunities for more rigorous approaches to exploit information sources of the type of the VOT. The complexity of the dependencies on various parameters like gender and phonetic context will therefore also be described experimentally.

Apart from applications in ASR, the current automatic VOT estimator can also be of interest in speech analysis, phonetics and speech pathology.

Acoustic information relevant to the identification of plosive sounds has been studied in the literature [4, 5, 6, 7]. Plosive consonants are produced by first forming a complete closure in the vocal tract via a constriction at the place of articulation, during which there is either silence or a low-frequency hum (called voicebar / prevoicing). The vocal tract is then opened, suddenly releasing the pressure built up behind the constriction. This opening of the vocal tract's airway is manifested acoustically by a transient and/or a short-duration noise burst. The duration of the interval between the release of the plosive and the beginning of voicing in the vowel is called the voice onset time or VOT. During this interval there is silence and/or noise caused by the release and/or aspiration noise. The VOT is one of the many acoustic cues for distinguishing plosives. The acoustic cues relevant to the articulation of a plosive can be related to manner (plosive, nasal, ...), place (bilabial, alveolar, velar, ...) and voicing (voiced, voiceless). A comprehensive discussion of these cues can be found in chapter 5 of [8] and we limit ourselves to an enumeration here. The *manner cues* for plosives include the presence of the silent region in the stop gap (obstruction phase), the rapid formant transitions and particularly a low locus frequency for the first formant F1, sudden energy change, release burst and aspiration. The *place cues* for plosives include the burst centre frequency (i.e. the main spectral peak of the turbulence occurring at the release), the locus frequency for the second and third formant transitions and the VOT. The *voicing cues* for plosives include the VOT, the presence of aspiration, the presence of an audible F1 transition, the intensity of the burst and the duration of the preceding vowel.

In this paper, we describe a VOT estimation algorithm using a high resolution signal analysis method which will better preserve timing information than MFCCs can. The next section is devoted to this signal representation, the reassigned time-frequency representation (RTFR). This representation allows to locate well-separated impulses, cosines and chirps in time and in frequency. Because speech can to some extent be seen as a sum of such signals, we advocate the use of this representation for our current task. In section 2.3, the VOT characteristics are highlighted. A VOT estimation algorithm starts with indentifying segments of speech that potentially contain a plosive sound. We therefore describe our plosive data sets in section 2.4 and move on to section 2.5 where the actual algorithm that computes the VOT feature from the RTFR is described. Although the VOT has already been studied extensively, there are not many algorithms described to *automatically* extract this feature. Related work can be found in [9, 10, 11, 12, 13]. However, to our knowledge this is the first time that the RTFR has been used to reliably extract the VOT feature. The performance of our algorithm is evaluated in section 2.6.1, while section 2.6.3 illustrates the modelling complexity as well as the usefulness of our automatic VOT extraction algorithm for phonetic studies by measuring some statistics of the VOT feature on the TIMIT database. Finally, in section 2.6.4 a rescored approach shows a modest improvement in speech recognition accuracy using VOT. Conclusions can be found in section 2.7.

2.2 Spectral reassignment

Time-frequency reassignment [14, 15, 16] offers an interesting solution for analysing transient signals such as plosives. The corresponding reassigned time-frequency representation (RTFR) has an increased sharpness of localisation of the signal components without sacrificing the frequency resolution. The RTFR is obtained by moving the spectral density value away from the point in the time-frequency plane where it was computed. The spectral density is reallocated from the geometric centre of the spectral analysis kernel function to the centre of gravity of the energy distribution. Though this principle can be applied to a multitude of time-frequency representations, here it is applied to the short time Fourier transform (STFT). Let $H(t, \omega)$, $D(t, \omega)$ and $T(t, \omega)$ denote the STFT of the signal obtained with the window function $h(t)$, the derivative of $h(t)$ and its time-weighted version $th(t)$ respectively and let $\Re(X)$ and $\Im(X)$ be the real and imaginary parts of X , then the energy at (t, ω) is reassigned to:

$$\hat{t} = t + \Re\left(\frac{T(t, \omega)}{H(t, \omega)}\right)$$

$$\hat{\omega} = \omega - \Im\left(\frac{D(t, \omega)}{H(t, \omega)}\right)$$

In practical implementations, the time-frequency plane is overlaid with a grid and reassigned energy is accumulated per cell.

In case the signal is a single cosine, linear chirp or Dirac impulse, the localisation in time and frequency is perfect. For instance, for a Dirac impulse $\delta(t - t_0)$ all energy will be reassigned to t_0 . When applied to speech with a sufficiently short analysis window, the RTFR clearly shows vertical (i.e. well-localized in time) lines for plosive bursts as well as for energy releases by the vocal folds. This property will make the construction of detectors for the burst onset of a plosive and for the subsequent start of periodicity (if any) fairly easy, as will be shown below. We have experimented with the multi-taper version of the RTFR [17], but a single window seemed to provide sufficient detail of the plosives to reliably reveal the acoustic events of interest, while it is computationally less demanding. Given the impulsive nature of the acoustic events we are trying to characterize, we opt for a Hamming window of length 8 ms, shifted by 0.625 ms per analysis frame. This corresponds to 128 and 10 samples respectively at a sampling frequency of 16 kHz which is adopted throughout this paper. Compared to the typical window lengths of 20 to 30 ms with a frame advance of 10 ms which are mostly used in speech recognition, our signal analysis offers a higher resolution in time. We used 256 equally spaced frequency bins for reassignment, a choice which is not critical given the wideband nature of the variables upon which the detection of the burst and the voicing onset will be based.

Figure 2.1 shows an example of the RTFR for a voiceless plosive (/t/) segment (followed by the vowel /ih/ as in "pit"), taken from the TIMIT database. The burst and onset of voicing as detected by the algorithm described in this paper are shown with arrows at the top. In this example, the burst of the /t/ is located at 15 ms, while the voicing starts at 87 ms, such that the VOT has a value of 62 ms. For comparison, we also show the original STFT from which the RTFR is computed in figure 2.2. Clearly, both the dental burst and the effects of glottal activity are better localized in time in the RTFR.

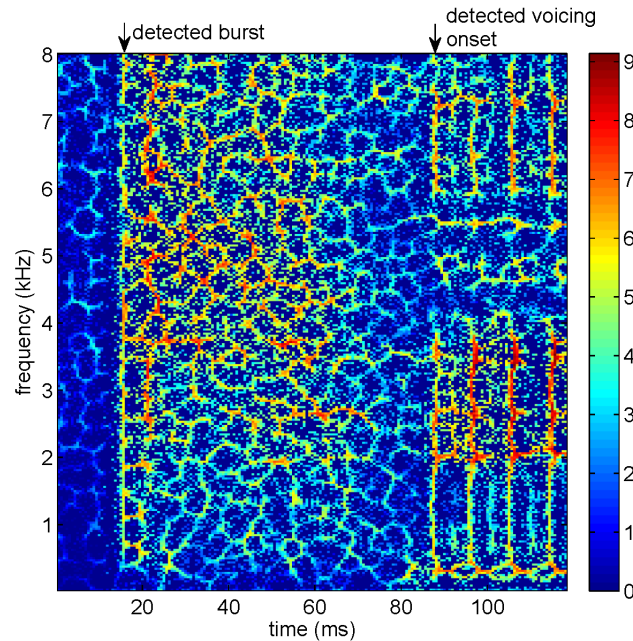


Figure 2.1: Reassigned time-frequency representation of a /t/ segment followed by /ih/. Color encode the logarithm of the energy.

2.3 Properties of the Voice Onset Time

On average, the VOT of voiceless plosives is larger than the VOT of voiced plosives, and the VOT increases from a bilabial to an alveolar and to a velar stricture. Hence, on average we have :

$$\begin{aligned} \text{VOT}(/bdg/) &< \text{VOT}(/ptk/) \\ \text{VOT}(/b/) &< \text{VOT}(/d/) < \text{VOT}(/g/) \\ \text{VOT}(/p/) &< \text{VOT}(/t/) < \text{VOT}(/k/) \end{aligned}$$

From the literature, we know that the VOT is influenced by several factors: the left and right context of the plosive, the position within the word, the lexical stress, speaker gender, speaking rate, the language, fundamental frequency F_0 of the vowel, ... For instance, there are notable differences in voicing across languages: Spanish has negative VOTs for the voiced plosives, while the VOTs of English are mostly positive. Women produce longer VOT values for voiceless stops than men [5]. Also, the VOT of children slightly changes with their age. When the plosive is followed by the vowel /i/, the mean VOT is larger than when it is followed by vowel /a/ [5]. An increase of the speaking rate causes a decrease of the VOT of voiceless plosives. Voiceless stops produced at a high fundamental frequency display shorter VOTs than those at low or mid F_0 's [6]. In addition, voiceless stops tend to display shorter VOTs and voiced stops display increased VOTs during conversational speech and reading, compared with single words.

Because of these effects, VOT distributions tend to overlap. Hence, the relation between the VOT value and plosive identity or even its voicing is not straightforward. Many studies try to circumvent this overlap by only considering plosives that are uttered in a constrained way, e.g. single words with a plosive in syllable initial pre-stressed position. In this way, the variability of the VOT within one class of plosives becomes smaller. In section 2.6, it will be shown that statistical models of the VOT

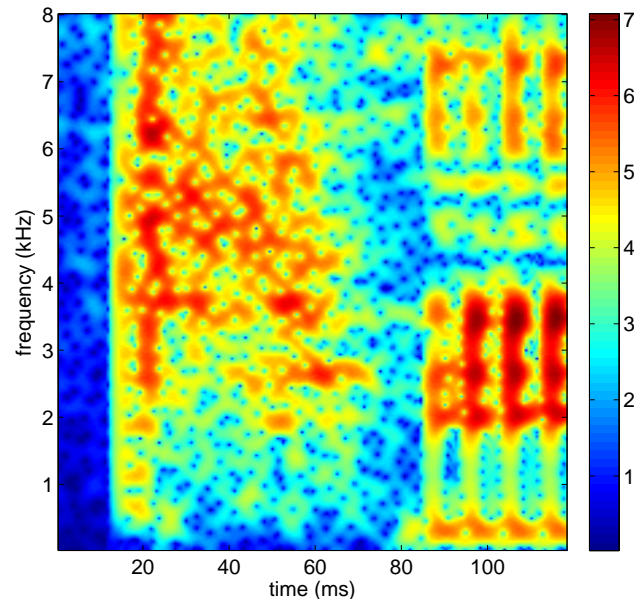


Figure 2.2: *STFT representation of the /t/ segment from figure 2.1. Color encode the logarithm of the energy.*

are more precise when they are conditioned on the phonetic context. If these models are to be used for accuracy gains in ASR as in section 2.6.4, the context can be assumed available (although not with 100% accuracy) from a first recognition pass when evaluating the estimated VOT. By using this knowledge, the overlap of the distributions can also be reduced to some extent.

2.4 Data sets

Experiments are conducted on the TIMIT database [18] since it contains manually verified phonetic transcriptions. It contains English read speech at office recording quality, uttered by native adults selected from eight dialect regions in the USA and sampled at a sampling frequency of 16 kHz. Though the algorithm may also apply to other plosives and affricates, this study focuses on the six plosives /p/, /t/, /k/, /b/, /d/ and /g/.

To study the quality of the VOT estimation algorithm that will be specified in (section 2.5), we adopt four data sets that are referred to as "forced", "manual", "free" and "test". Each of these sets contains a collection of segments of speech in which we expect to find one of the six plosives. Depending on the data set, the segment identity as well as its boundaries are generated in different ways as described below. The number of speech segments for each plosive is given in table 2.1.

2.4.1 The "forced" data set

The "forced" data set is relevant for phonetic studies, for automated studies of the parameters affecting the VOT or for automated pronunciation scoring in (foreign) language learning. In these settings, speech segments can be found in which one of the plosives under study is present and our task is to

estimate the VOT. The segment boundaries are obtained from a forced alignment with an HMM-based speech recogniser using the manually verified phonetic transcriptions available in the TIMIT database. Hence, we rely on information that is normally not available in an automatic speech recognition system. All occurrences of the six plosives from the 3696 phonetically rich "si" and "sx" training utterances originating from 462 different speakers in the TIMIT database are included in the "forced" data set, irrespective of the left and right phonetic context.

The acoustic models used for segmentation are context independent HMMs with 2 to 4 states per phone trained on an independent data set. In total, there are 141 GMMs sharing 5550 Gaussians with diagonal covariance. The speech features are mel-scaled log-filterbank outputs that are linearly transformed with a decorrelating and diagonalizing transform [19]. Since these features are recalculated every 10 ms, this is also the segmentation resolution. Voiced plosives and voiced affricates share a common 2-state HMM for the closure. The voiceless plosives and affricates also share their closure model. By including separate models for the phone components of plosives, the HMM will produce separate segments for the closure and the burst. The segment boundaries that are associated with the plosive are those of the burst only. The reason for this choice is that the segment boundaries generated by the HMM will serve as a fallback in case we fail to detect the burst or the onset of voicing, while the duration of the burst segment can be seen as a measurement of the VOT.

2.4.2 The "free" data set

In a fully automatic VOT extraction setting, a forced alignment is not possible due to the lack of a unique transcription hypothesis. Therefore, in the second data set, plosive segment candidates are generated by a phonetic automatic speech recogniser as described in [20] applied to the same utterances used in the "forced" data set. The HMMs described in section 2.4.1 are used to find the best matching phonetic transcription using a phone-level bigram language model with Witten-Bell smoothing [21]. Any segment automatically labeled as the burst of one of the six plosives under study was included in the set, irrespective of the detected phone or phone component on the left and on the right.

2.4.3 The "manual" data set

The performance of the algorithm will be evaluated by comparing the automatic VOT estimates with values derived by an expert. To this end, a subset of the plosive speech segments was selected from the "forced" set as follows. Cycling through all 16 gender/dialect combinations, we randomly drew a speaker from that gender/dialect combination and subsequently we randomly drew a recording (sample file) from that speaker. For any of the six plosives for which we collected less than 130 examples so far, the expert manually estimated the VOT of all occurrences in the recording by inspection of waveforms and spectrograms centered around the automatically generated segment boundaries, marking the burst onset time and the start of voicing and finally storing the time difference. In total 268 different recording files from the TIMIT database were used. All plosive segments that were not followed by a voiced sound or for which the manual annotator could not detect a burst or the start of voicing were removed. There is no constraint on the left phonetic context. Table 2.1 shows the exact number of examples thus retained in the "manual" data set.

Table 2.1: *Number of speech segments in each of the data sets.*

	forced	free	manual	test
/b/	2181	2012	115	754
/d/	2432	2222	76	728
/g/	1191	977	98	386
/p/	2588	2749	111	821
/t/	3948	4052	92	1180
/k/	3794	3968	90	1039
total	16134	15980	582	4908

2.4.4 The "test" set

This set is constructed exactly like the "forced" data set, except that the sentences are taken from the TIMIT test set ("extended" set without the "core" set), a total of 1152 sentences from 144 speakers.

2.5 The VOT estimation algorithm

The actual estimation of the VOT requires that the burst onset and the start of periodicity are determined for each plosive segment. In this section, we will describe how the segment boundaries are obtained and how both events are detected. The process is illustrated in figure 2.4. The estimated VOT is then the elapsed time between the estimated burst onset time and the estimated start of periodicity.

2.5.1 Detection of plosive segments

The first step in the algorithm consists of finding segments in the speech signal that could contain a plosive. Such segments could be found using dedicated detectors, as is shown in the research on automatic extraction of phonological features. In [22] and [23], detectors are described that exhibit sufficient accuracy to generate candidate plosive segments. The method used for generating plosive segment candidates is important to the performance of the algorithm for three reasons. First, segments may be missed or overgenerated, leading to unrecoverable errors. Second, we will search for burst and voicing within boundaries derived from the proposed segment. Errors in the segment boundaries may cause to wrongly identify an acoustic event as the burst or voicing onset. Third, in case either burst or voicing cannot be detected automatically, fallback estimates of their time of occurrence are derived from the proposed segment boundaries.

In the current work, we have opted for a HMM-based automatic speech recogniser to generate plosive segment candidates, as explained in section 2.4. Depending on the application of the VOT estimate, it may or may not be reasonable to assume that a phonetic transcription of the speech around the plosive is available. We therefore defined the "forced" and "free" data sets in which plosive segments are generated with or without phonetic knowledge of the test utterance. In both sets, the algorithm will start looking for the burst 2.5 ms or 4 frames prior to the burst segment start found by the recogniser. Starting earlier would increase the risk of misdetecting energy bursts from the previous phone as belonging to the plosive. Starting later would increase the risk of missing the burst. The end of the segment is extended by 10 ms or 16 frames to the future. Extension of the segment end to the right just means more pitch cycles will be included and is harmless to the algorithm. The value of 10 ms is a compromise such that at least one glottal closure will be seen in most cases, while avoiding

unreasonably high VOT estimates in case some initial glottis vibration cycles are not detected. Notice that even if the successor segment was manually or especially automatically labeled as a vowel, this does not *guarantee* that glottal activity will be detected.

In the discussion below, we will refer to *extended* segments to refer to the plosive segment starting 2.5 ms before and ending 10 ms after the segment determined by the speech recogniser.

2.5.2 Burst onset detection

Figure 2.1 shows that the onset of the release phase gives rise to a sudden increase of the amplitude over the whole frequency range.

To limit the influence of the high-amplitude pitch pulses which also have a strong low-frequency component, only the frequency range 3.2-8 kHz is retained for burst detection. The corresponding frequency bins in the RTFR power are summed to form the "burst power" $p(n)$ estimate for frame n . Then, the first local maximum that is sufficiently strong and ramps up sufficiently sharply is identified as the burst onset. The condition is asymmetric because $p(n)$ can stay high during the release interval after the burst. In formulae, frame n is retained as a possible burst location if it satisfies all of the following conditions $p(n) > p(n - j)$, for $j = -1, 1$ and 2 (local maximum), $p(n) - p(n - i) > p_m(n)$ for $i = 2 \dots 5$ (sufficiently sharp and strong peak), where $p_m(n)$ is a measure that relates to the average signal energy so the criteria are invariant to scaling of the signal. In our experiments, $p_m(n)$ is taken to be the mean of $p(n)$ over 150 plosive frames.

If the automatic algorithm does not find a local maximum, the start of the (unextended) segment is marked as the burst onset. This may happen because the burst is simply missing (by construction, this will not happen in the "manual" data set) or because it is too weak. The resulting estimate is less accurate: measured over all plosives of the "manual" data set, the square root of the mean square estimation error is 12.6 ms if a burst was detected, while it increases to 22.6 ms if a burst could not be detected.

2.5.3 Start of periodicity

As can be seen from the RTFR in figure 2.1, the periodicity of the signal gives rise to vertical lines of high amplitude with valleys in between. The distance between these lines is determined by the pitch period. This periodic structure is mainly present in the lower part of the frequency range.

To obtain a robust estimate of the start of voicing, only the frequency range 0-4 kHz is retained. At a sampling frequency of 16 kHz as used in this work, this comes down to keeping only the lower half of the RTFR. Then, a short term autocorrelation is computed by multiplying every RTFR frame (for every 0.625 ms frame advance) with a weighted version of the frames at lags 1 to 40 and summing these values over the lag index and over the retained frequency bins. The weighting function (figure 2.3) is given by the difference of two decaying exponential functions and has a large value in the adult pitch period range of 5 to 20 frames, corresponding to a pitch period between 3.1 ms and 12.5 ms or a pitch frequency of 80 Hz down to 320 Hz. An asymmetric weighting function is chosen because we want to extract the *start* of periodicity. The result is normalised with the total energy in the frames under the autocorrelation window over the whole frequency range (0-8 kHz).

The autocorrelation function obtained in this way will exhibit a large value at times where there is a substantial amount of energy that is periodically repeated within the analysis frame, i.e. at the time instants for which a pitch pulse is present in the RTFR. To be marked as a local maximum, the following conditions have to be met : its value has to be larger than the value of its direct neighbours,

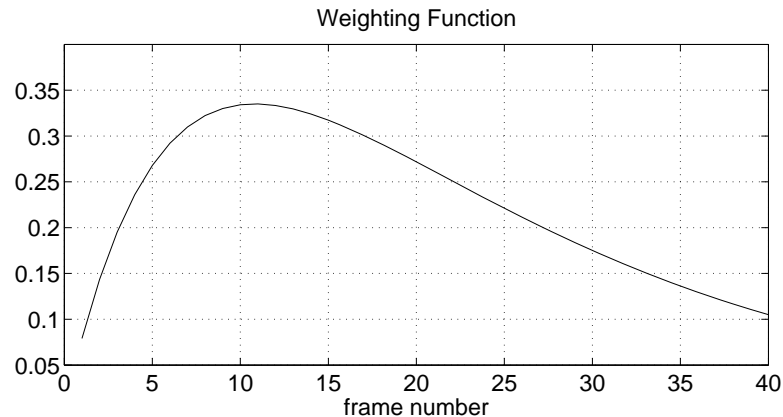


Figure 2.3: *Weighting function of the periodicity detector.*

and it has to exceed the value of its neighbours at distances of ± 2 , ± 3 and ± 4 frames with an increasing threshold to assure that the selected peaks are at least 5 frames (or the minimum pitch period) on either side from their neighbours and at least 0.03 in height, a value which was determined from visual inspection on the "forced" data set (excluding the "manual" set).

With this scheme, some of the bursts will also be marked as pitch pulses. Moreover, a velar stricture can have multiple bursts that should not be confused with pitch pulses. To avoid selecting the burst as the start of voicing, an additional constraint is imposed. A local maximum has to be within the maximal pitch period (20 frames or 12.5ms) from the *next* local maximum (or from the end of the extended segment). For low-pitched voices, the wrong starting point of voicing can still be selected if some pitch pulses are not detected. However, the risk of selecting the burst onset is strongly reduced, especially if multiple bursts are present.

If the algorithm cannot detect voicing within the extended segment, the end of the unextended segment is marked as the start of voicing, i.e. we fall back to the HMM's decision of the start of the next phone. This is a reasonable choice for English, where VOTs are mostly positive, but for other languages, voicing may already start in the closure interval. On the "manual" data set, we measure a square root of the mean square error of 12.2 ms if voicing was detected, while it increases to 17.8 ms if voicing could not be detected within the extended segment. Not surprisingly, the HMM does a better job at detecting the start of the next vowel than it does at detecting the burst.

2.5.4 Discussion

The proposed peak picking algorithms are surely not the only possible approaches to detecting the burst and voice onset events in RTFRs. The advantage of the RTFR is that the peaks are clear and sharp, which motivates the high time resolution of 0.625 ms used in our proposed algorithm. Often, both the burst and the glottal closures can be located to a single frame. Decreasing the frame rate might make the algorithm computationally more efficient, but would make the peak picking more error prone. In any case, even at pitch periods down to about 3 ms, sampling needs to be fast enough to resolve the pitch peaks. Similarly, the burst onset may exhibit multiple clicks which should not be merged into a single broad peak of $p(n)$ if the same peak detection criteria are maintained.

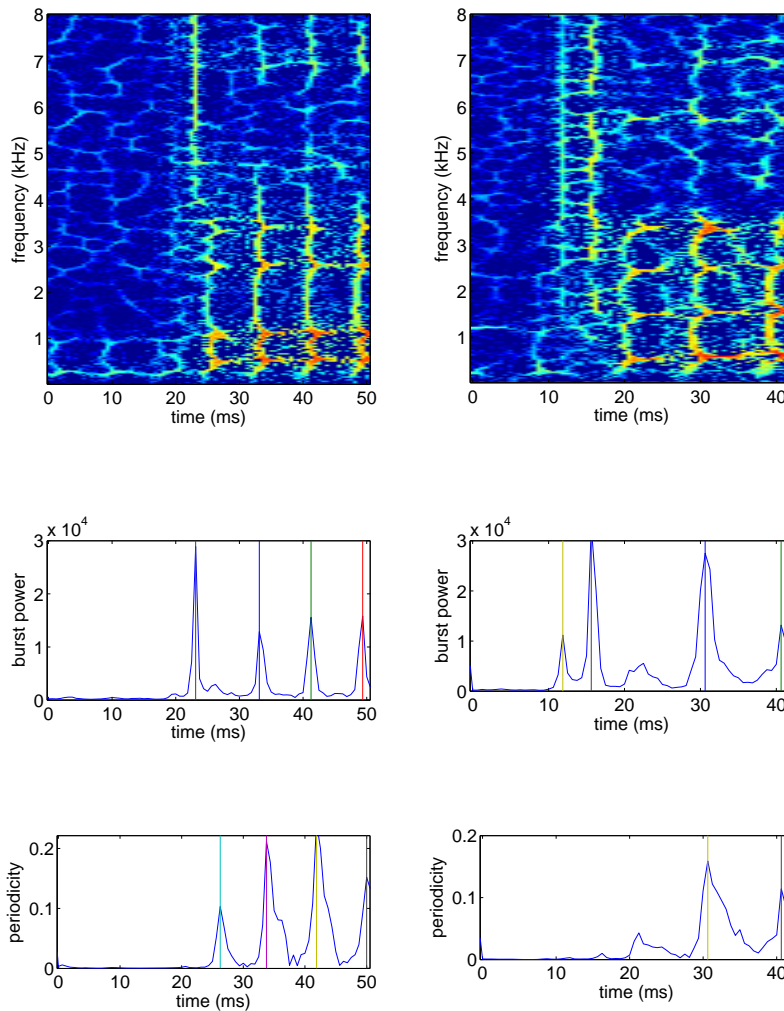


Figure 2.4: *Left: illustration of the peak picking on a /b/ segment with a right context /aa/ (from "flat bottom") taken from the "free" data set. From top to bottom: RTFR, burst detection and periodicity detection. The peaks that satisfy the selection criteria are marked with vertical lines. Right: /b/ segment (from the word "thereby") with erroneous detection of the start of voicing.*

2.6 Experiments

2.6.1 Algorithm performance for phonetic studies

The VOT was estimated for the complete "forced" data set by means of the automatic algorithm of section 2.5. Since the "manual" data set is a subset of the "forced" set, it is possible to compare the manual and automatic VOT estimates on this subset. Figure 2.5 shows the cumulative distribution of the absolute difference between the manually and the automatically extracted VOT estimates. On average, the difference is smaller than 10 ms in 76.1% of the plosive segments, smaller than 20 ms for 91.4% of the plosive segments, and smaller than 30 ms for 96.2% of the plosive segments. The average deviation from the manually assigned VOT is the largest for /d/ and decreases from /d/ to /k/, /g/, /t/, /p/ and /b/.

Table 2.2 gives an indication of the bias of the algorithm. For each plosive, it contains the average of the manually and of the automatically extracted VOTs on the "manual" data set. The resulting bias is calculated as the difference of both means and the uncertainty on this estimate is given as its standard deviation assuming independent bias measurements. There is an overall bias of 2.9 ms, which is even statistically detectable on most individual plosives. To show that the bias is mainly due to the fallback in case either burst or voicing onset cannot be detected automatically, the right side of the table gives the same statistics measured only on those segments from the "manual" data set for which the algorithm was able to detect both events. The overall bias is now down to 0.9 ms and mostly realized on /d/. A further analysis would need to question the human annotation as well as the peak selection criteria. Phenomena as illustrated in the right pane of figure 2.4 are likely to play a role here.

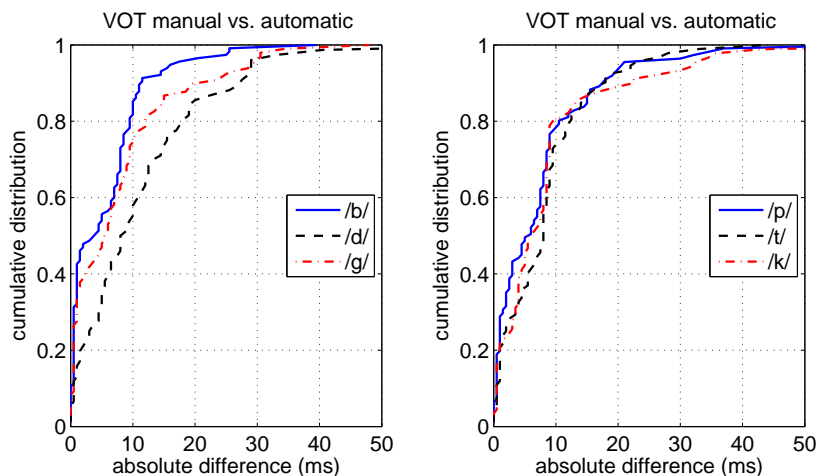


Figure 2.5: Absolute difference between the manually and the automatically extracted voice onset time.

2.6.2 Algorithm performance for automatic speech recognition

While the above accuracy analysis is relevant for e.g. phonetic studies, where segment boundaries can be generated based on a manually produced phonetic transcription, its validity can be questioned in a fully automatic setting, where the goal of VOT estimation could be to improve speech recognition

Table 2.2: Comparison between the average manually and automatically extracted VOT for each plosive. Left: all plosive segments of the "manual" data set. Right: only the plosive segments for which both burst and voicing onset could be detected automatically.

	VOT (ms) all segments				VOT (ms) without fallback			
	manual	autom	bias	stdev	manual	autom	bias	stdev
/b/	7.7	9.8	2.1	0.9	7.9	8.8	0.9	0.8
/d/	8.5	16.1	7.7	1.9	8.2	13.5	5.2	2.0
/g/	21.8	22.7	0.9	1.1	21.7	21.7	0.0	1.1
/p/	39.4	44.1	4.6	1.1	38.5	40.4	1.9	1.2
/t/	50.9	51.4	0.6	1.2	50.2	48.9	-1.3	1.3
/k/	54.3	56.4	2.1	1.7	56.2	55.2	-1.1	2.0
avg	30.3	33.1	2.9	0.5	28.8	29.7	0.9	0.5

accuracy on plosives. Therefore, in the second study, the absolute difference between manual and automatic estimates is analysed on the "free" data set. However, an automatic phone recogniser can mislabel plosive segments, insert or omit them, or generate different segment boundaries. We related the plosive segments from the "free" data set with one from the "forced" data sets by selecting the "forced" plosive segment with the largest overlap in time. For 9.2% of the segments, there was no overlap. Only 0.04% of "free" segments overlapped with more than one "forced" segment, in which case we took the "forced" plosive with the largest overlap in time. Notice that it may well be that the phone identity (among the set of six considered) is different in both sets, corresponding to the mislabelings by the recogniser that we are trying to correct. In this analysis, the manual phonemic labelings provided the TIMIT database are assumed to be correct.

With this procedure, 566 plosive segments from the "free" set could be linked with a segment from the "manual", which allows to recompute the cumulative distribution of the absolute difference between manual and fully automatic VOT estimates. The percentiles for 10 ms, 20 ms and 30 ms deviation now become 72.6%, 87.8% and 93.8% respectively (instead of 76.1%, 91.4% and 96.2%). Hence, the main source of estimation error is not caused by the automatic generation of segment boundaries. Also notice that only 16 (= 582 - 566) out of 582 plosive segments from the "manual" set could not be found automatically, which is far less than 53 (9.2 % of 582). Hence, the HMM-based plosive detector performs a lot better on plosives for which the human annotator found a burst and that are followed by a voiced sound.

2.6.3 Estimated VOTs

With this automatic algorithm, we can investigate to which extent factors such as gender and phonetic context could be taken into account in statistical models. In this study, we focus on the voicing dimension, rather than place of articulation.

First, we measure the effect of gender. The second column of table 2.3 shows the VOT obtained on the "forced" data set for each of the plosives, averaged over all speakers and all contexts. These values confirm the inequalities of section 2.3. Columns 3 and 4 contain the VOT values averaged over all contexts but including only the male, or only the female speakers, respectively. On our database, the VOTs of plosives uttered by women are on average 12% longer than that of men. For /p t k/, this is in line with [5], but the latter article did not mention the same effect for /b d g/. Notice that

Table 2.3: VOT estimate [ms] for each plosive class, averaged over all contexts in the "forced" data set. Mean value for all speakers, only male or only female speakers. Columns 5-7 indicate the corresponding number of segments.

	VOT [ms]			# segments		
	m + f	m	f	m + f	m	f
/b/	11.8	11.3	13.0	2181	1522	659
/d/	18.6	17.7	20.5	2432	1681	751
/g/	21.8	20.7	24.0	1191	800	391
/p/	40.8	39.0	45.0	2588	1798	790
/t/	43.6	41.8	48.1	3948	2791	1157
/k/	48.0	47.1	50.3	3794	2686	1108

the gender-independent averages differ from those of table 2.2 because the phonetic context of the plosives differs, as explained in section 2.4.

The effect of the right context can be found in figure 2.6, which presents the VOT means together with the standard deviations without any right context imposed or when it is followed by a vowel /ih/ (as in "bit"), /aa/ (as in "box") or /eh/ (as in "bet"). There is no constraint on the left context. In total, there are between 68 and 253 examples of each right-context dependent plosive in the database when pooling over all speakers. If the phonetic context is constrained, the overlap of the VOT distributions usually decreases. For instance, the error bars of /k eh/ and /g eh/ do not overlap, while the error bars for the context independent /k/ and /g/ do. The same can be said about /p aa/ and /b aa/ versus /p/ and /b/. The longer average VOT for right context /ih/ than for context /aa/ is only observed for plosives /b d g t/.

Figure 2.7 shows histograms of the context dependent VOTs of plosives followed by the vowel /eh/, constructed on the "forced" data set. From this figure, the overlap of the distributions is clearly apparent. This overlap is even larger for the context independent histograms. This illustrates that the relation between the VOT value and the voicing cue of the plosive is not straightforward.

2.6.4 VOT as a feature for automatic speech recognition

Histograms like the one of figure 2.7 can be used in a likelihood ratio test to discriminate, for instance, along the voicing dimension. To this end, context dependent but gender independent histograms are built with 23 uniformly spaced bins 5 ms apart between -10 ms and +100 ms using the "forced" data set. Let $N(V, l, p, r)$ be the number of times the estimated VOT falls in bin V for plosive p with left context l and right context r . Overall, 1298 different phone/plosive/phone combinations are observed. Many of these histograms have little data, so a multi-stage backoff scheme is applied to histograms having less than 40 counts, i.e. if

$$\sum_V N(V, l, p, r) < 40$$

First the left context is generalised to one of 12 broad phonetic classes, then the right context is generalized, then the left context is disregarded and finally the right context is disregarded. The back-off steps are terminated as soon as at least 40 counts are observed in the histogram with the generalized context. We will call the thus obtained generalized left and right context \tilde{l} and \tilde{r} respectively.

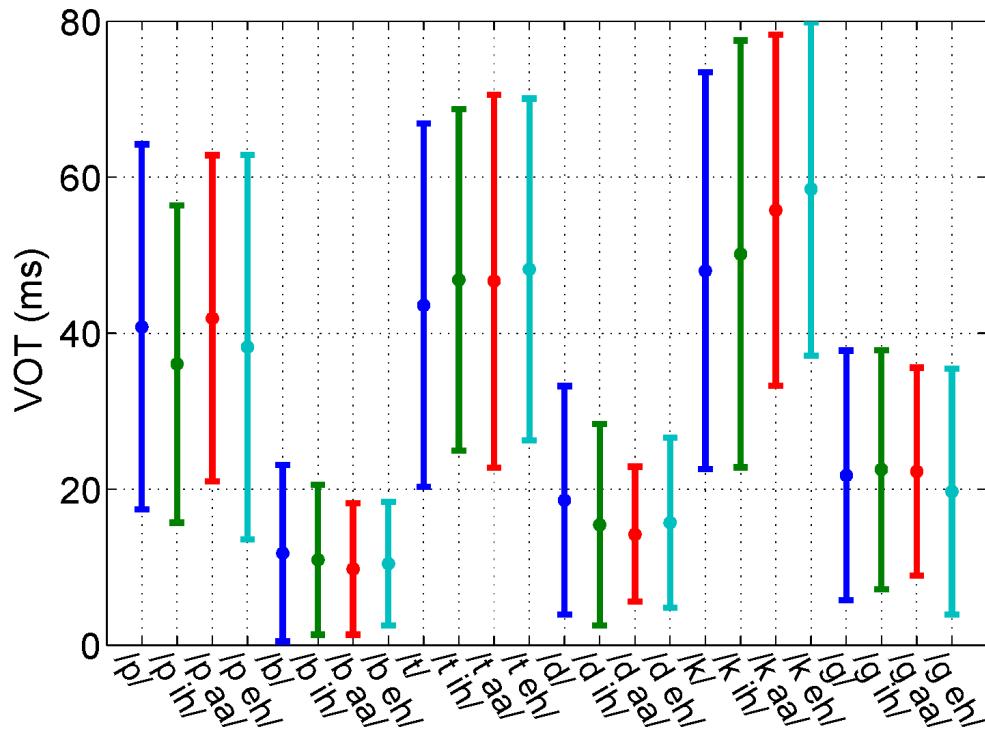


Figure 2.6: Mean VOT for plosives /p b t d k g/ by context (context independent, right context /ih/, /aa/, /eh/). The left context is always unconstrained. Error bars indicate +/- one standard deviation. Measured on the "forced" data set.

Figure 2.8 shows the logarithm (to base 10) of the likelihood ratio versus the estimated VOT value for the "test" data set. This set contains data that was not used during the construction of the histograms, while the ground truth about plosive identity and its context is known from the manual labeling provided in the TIMIT database. So let $P(V|l, p, r)$ be the probability that the estimated VOT falls in bin V for plosive p as measured on its histogram, and let $P(V|l, \bar{p}, r)$ be the probability read from the histogram for the plosive with opposite voicing. The log-likelihood ratio is then

$$\log_{10} \left(\frac{P(V|l, p, r) + \varepsilon}{P(V|l, \bar{p}, r) + \varepsilon} \right)$$

where

$$P(V|l, p, r) = \frac{N(V, \tilde{l}, p, \tilde{r})}{\sum_V N(V, \tilde{l}, p, \tilde{r})}$$

and ε is a small constant to cope with zero probability estimates and was set to 10^{-3} in our experiments. The left panes show the log-likelihood ratio on the voiceless data and assuming the voiceless sound (p is /p/, /t/ or /k/ and \bar{p} is /b/, /d/ or /g/ respectively), while the right panes show the log of the reciprocal on the voiced data (i.e. assuming p is a voiced sound). Figure 2.8 illustrates that large (small) VOTs for voiceless (voiced) sounds indeed lead to positive log-likelihood ratios, but that negative log-ratios can occur. That the choice of ε is not a critical one is also apparent from

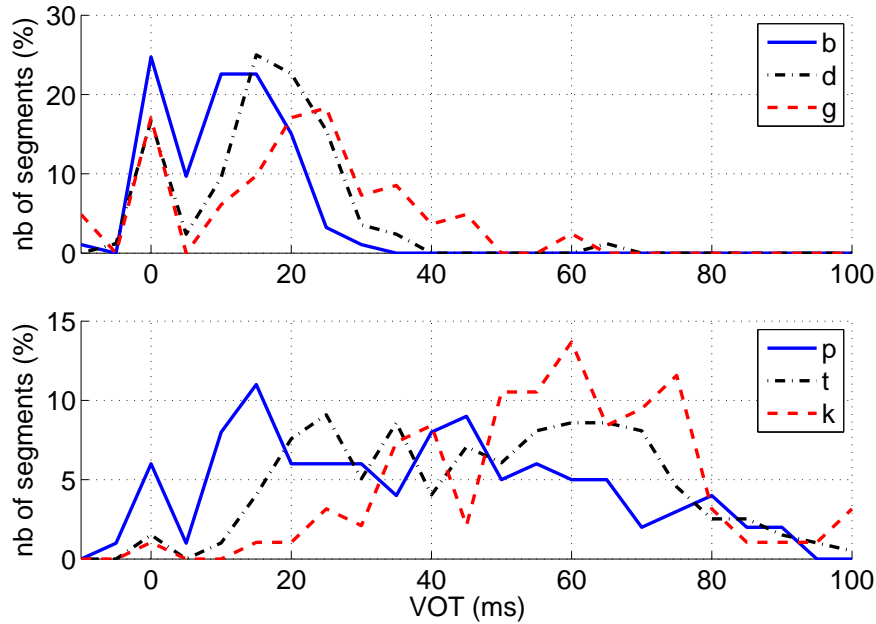


Figure 2.7: Normalised histogram of VOT estimates on the "forced" data set for plosives /b d g/ and /p t k/ followed by vowel /eh/, without constraint on the left context.

these scatter plots. Its side-effect is to limit extreme values of the log-likelihood ratio, an effect that is mostly observed on the positive side.

In an attempt to improve the phone recognition rate by exploiting the VOT as a feature, phone lattices were generated on the TIMIT test data as described in [20]. These are the same sentences as used in the "test" data set, but now the lattice will include more plosive candidates. The best path through the lattice will generate the phone segmentation of the "test" data set. In formula 2.1 (see below), the likelihood $L(A)$ of each plosive arc A in the lattice is then linearly combined with the log-likelihood ratio of it being correct versus its variant with opposite voicing being correct. There is, however, a difference with the above. When dealing with the "test" data set, the left and right phonetic context are unique. In a lattice, multiple arcs may arrive in the starting node of A and multiple arcs may leave from its ending node, so the left and right phonetic context are not unique. We denote the set of phone labels of arcs ending (starting) in the starting (ending) node of arc A with \mathcal{L} (\mathcal{R}) and sum the statistics over all contexts of A allowed by the lattice:

$$P(V|\mathcal{L}, p, \mathcal{R}) = \frac{\sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}} N(V, \tilde{l}, p, \tilde{r})}{\sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}} \sum_V N(V, \tilde{l}, p, \tilde{r})}$$

The corrected acoustic likelihood of a lattice arc A becomes:

$$L(A) + \alpha \log_{10} \left(\frac{P(V|\mathcal{L}, p, \mathcal{R}) + \varepsilon}{P(V|\mathcal{L}, \bar{p}, \mathcal{R}) + \varepsilon} \right) \quad (2.1)$$

Linear combination of log-likelihoods of different information sources was examined in [24]. The single free parameter α we introduced was tuned on the "forced" data set, which is independent of the

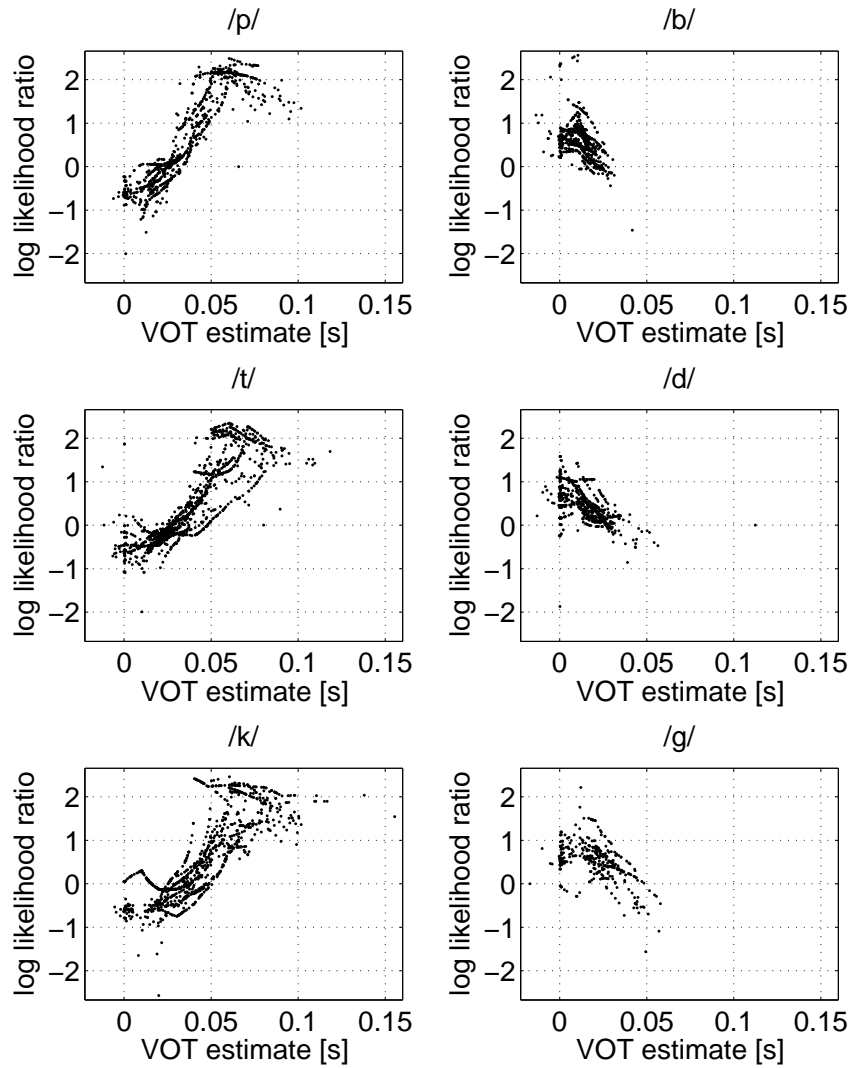


Figure 2.8: *Logarithm of the likelihood ratio versus the automatically calculated VOT value, measured on the "test" data set.*

"test" data set. This procedure reduced the phone error rate from 26.70% to 26.53% on the TIMIT test set. Hence, we observe that the VOT feature has contributed only very little to error rate improvement. This is not surprising, since we observe in figure 2.8 that the log-likelihood ratio can become negative for valid utterances of the plosive. On the other hand we have to realize that we attempt to correct only the plosive hypotheses generated by the HMM system, and this only along the voicing dimension. We can find the best obtainable error rate by correcting the voicing of the plosives in the first best path through the phone lattice using the reference transcription. This yields an error rate floor of 25.85%. Hence, we have obtained $(26.7 - 26.53)/(26.7 - 25.85) = 20\%$ of the performance gain that would be achievable using an ideal voicing detector. In absolute numbers, the VOT-based likelihood ratio test corrected 80 out of 1853 plosive errors and hence the improvement is statistically significant. The gain shows that the VOT estimate does contain information that the HMM is not able to exploit. Apart from the overlap in the distributions of the VOT, the performance in this particular implementation is also limited by the pruning in the phone lattice. Each plosive hypothesis (arc) is rescored, but this can only lead to a change in decision if the hypothesis with opposite voicing is also in the lattice (and receives a better combined score). Hence, if the alternate, correct hypothesis was not included in the lattice because of pruning, it cannot be recovered, even with an ideal voicing detector. Further performance improvements might also be obtained by combining the HMM and VOT likelihoods in a non linear way.

2.7 Conclusions

We described an algorithm to *automatically* extract the voice onset time. It operates on the reassigned time-frequency representation of the signal, which has an improved localisation of the relevant acoustic events. The algorithm performance was characterised for English plosives on the TIMIT database. The accuracy seems sufficient to reconstruct some of the findings of the literature on phonetics about the factors affecting VOT. Using a rescoreing approach, it was shown that the automatic VOT estimate does provide some additional information about the phone identity which is not exploited in state-of-the-art HMM-based ASR systems.

Bibliography

- [1] C.-H. Lee, M. A. Clements, S. Dusan, E. Fosler-Lussier, K. Johnson, B.-H. Juang, and L. Rabiner, "An overview on automatic speech attribute transcription (asat)," in *Proc. INTERSPEECH*, Antwerp, Belgium, Sep. 2007, pp. 1825–1829.
- [2] D. Seppi, D. Falavigna, G. Stemmer, and G. R., "Word duration modeling for word graph rescoring in lvcstr," in *Proc. INTERSPEECH*, Antwerp, Belgium, Sep. 2007, pp. 1805–1808.
- [3] J. Bilmes and C. Bartels, "Graphical model architectures for speech recognition," *IEEE Signal Processing Magazine*, vol. 22, no. 5, pp. 89–100, 2005.
- [4] S. O'Brien, "Spectral features of plosives in connected-speech signals," *International Journal on Man-Machine Studies*, vol. 38, pp. 97–127, 1993.
- [5] S. Whiteside, L. Henry, and R. Dobbin, "Sex differences in voice onset time: A developmental study of phonetic context effects in british english," *Journal of the Acoustical Society of America*, vol. 116, no. 2, pp. 1179–1183, 2004.

- [6] C. McCrea and R. Morris, "The effects of fundamental frequency level on voice onset time in normal adult male speakers," *Journal of Speech, Language, and Hearing Research*, vol. 48, pp. 1013–1024, 2005.
- [7] J. Jiang, M. Chen, and A. Alwan, "On the perception of voicing in syllable-initial plosives in noise," *Journal of the Acoustical Society of America*, vol. 119, no. 2, pp. 1092–1105, 2006.
- [8] G. J. Borden and K. S. Harris, *Speech Science Primer: Physiology, Acoustics, and Perception of Speech*, 2nd ed. Baltimore, U.S.A.: Williams & Wilkins, 1984.
- [9] C. Lefebvre and D. Zwierzynski, "The use of discriminant neural networks in the integration of acoustic cues for voicing into a continuous-word recognition system," in *Proc. International Conference on Spoken Language Processing*, Kobe, Japan, Nov. 1990, pp. 1073–1076.
- [10] P. Ramesh and P. Niyogi, "The voicing feature for stop consonants: Acoustic phonetic analyses and automatic speech recognition experiments," in *Proc. International Conference on Spoken Language Processing*, Sydney, Australia, Nov. 1998.
- [11] P. Niyogi and P. Ramesh, "Incorporating voice onset time to improve letter recognition accuracies," in *Proc. International Conference on Acoustics, Speech and Signal Processing*, Seattle, WA, U.S.A., May 1998, pp. 13–16.
- [12] K. Sonmez, M. Plauche, E. Shriberg, and H. Franco, "Consonant discrimination in elicited and spontaneous speech: a case for signal-adaptive front ends in ASR," in *Proc. International Conference on Spoken Language Processing*, Beijing, China, Oct. 2000.
- [13] A. Kazemzadeh, J. Tepperman, J. Silva, H. You, S. Lee, A. Alwan, and S. Narayanan, "Automatic detection of voice onset time contrasts for use in pronunciation assessment," in *Proc. International Conference on Spoken Language Processing*, Pittsburgh, PA, U.S.A., Sep. 2006.
- [14] F. Auger and P. Flandrin, "Improving the readability of time-frequency and time-scale representations by the reassignment method," *IEEE Transactions on Signal Processing*, vol. 43, no. 5, pp. 1068–1089, 1995.
- [15] F. Plante, G. Meyer, and W. Ainsworth, "Improvement of speech spectrogram accuracy by the method of reassignment," *IEEE Transactions on Speech and Audio Processing*, vol. 6, no. 3, pp. 282–286, 1998.
- [16] S. Hainsworth and M. Macleod, "Time-frequency reassignment: a review and analysis," Cambridge University Engineering Department, Tech. Rep. CUED/FINFENG/TR.459, 2003.
- [17] J. Xiao and P. Flandrin, "Multitaper time-frequency reassignment for nonstationary spectrum estimation and chirp enhancement," *IEEE Transactions on Signal Processing*, vol. 55, no. 6, pp. 2851–2860, 2007.
- [18] J. S. Garofolo and et al., "DARPA TIMIT acoustic-phonetic continuous speech corpus CDROM," 1993.
- [19] K. Demuynck, "Extracting, modelling and combining information in speech recognition," Ph.D. dissertation, K.U.Leuven, Belgium, Feb. 2001.

- [20] K. Demuynck, D. Van Compernelle, and H. Van hamme, “Robust phone lattice decoding,” in *Proc. International Conference on Spoken Language Processing*, Pittsburgh, U.S.A., Sep. 2006, pp. 1622–1625.
- [21] I. Witten and T. Bell, “The zero-frequency problem: Estimating the probabilities of novel events in adaptive text compression,” *IEEE Transactions on Information Theory*, vol. 37, no. 4, pp. 1085–1094, 1991.
- [22] S. King and P. Taylor, “Detection of phonological features in continuous speech using neural networks,” *Computer Speech and Language*, no. 14, pp. 333–353, 2000.
- [23] F. Stouten and J.-P. Martens, “Speech recognition with phonological features: Some issues to attend,” in *Proc. International Conference on Spoken Language Processing*, Pittsburgh, PA, U.S.A., Sep. 2006, pp. 357–360.
- [24] P. Beyerlein, “Discriminative model combination,” in *Proc. International Conference on Acoustics, Speech and Signal Processing*, Seattle, WA, U.S.A., May 1998, pp. 481–484.

Chapter 3

Using Prosody with NMF

G. Aimetti (Sheffield), H. Van hamme (KU Leuven)

3.1 Introduction

The aim of these experiments was to investigate whether the addition of prosodic cues, rhythm and pitch, would increase word detection accuracy for the nonnegative-matrix factorization (NMF) approach. It was hypothesised that the addition of prosodic cues to the input stream would increase the accuracy at a faster rate.

3.2 Rhythm Vectorisation

The first prosodic cue to be employed was rhythm. The *Rhythmogram* model of Todd and Brown [1] is used to detect peaks within the speech signal (Fig. 3.1), by looking over a range of different time constants in order to derive hierarchical structure of the onsets of individual events.

The amplitude envelope of the utterance is calculated which is then passed to a multi-scale Gaussian low-pass filter system. The rhythmogram output is derived by finding the peaks in the low-pass response or zero crossings of the 1st derivative. An example is provided in Fig. 3.1 (right). Time is displayed across the horizontal axis and the different time constants used are on the vertical axis.

The NMF technique requires the data to be decomposed and made available in the form of a data matrix. This implies that an additional stream of information must be encoded in terms of a sequence of vectors. To vectorise the rhythm for NMF we exploited the manner in which syllabic-type events converge to higher levels within the utterance. We achieved this by creating event strings across all channels with varying time constant. So, for each event we collect the time distance of the nearest peak in each consecutive time constant channel, where the root position is channel 1 (time constant = 1ms). Figure 3.2 shows the event strings within an example utterance.

The time distances are then quantized and labelled in terms of a value between 1 and 20 (see 3.2), with finer resolution for shorter distances. The result is a histogram in terms of the labels 1-20. The final stage of the rhythm vectorisation process is to create a co-occurrence matrix of the quantised peak distances within the utterance (Fig. 3.3). This stream can now be appended to the NMF input.

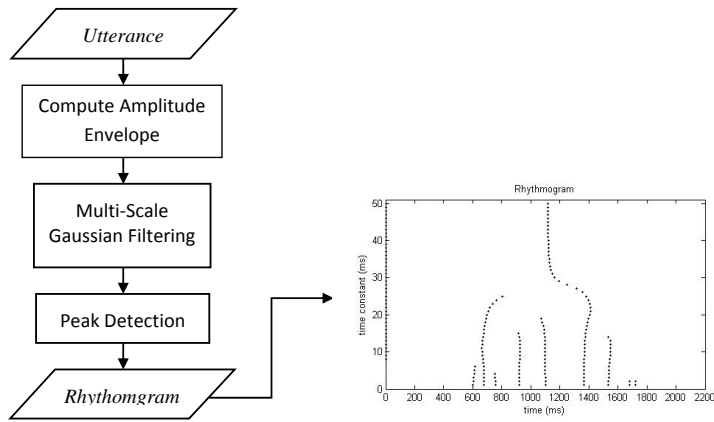


Figure 3.1: *Rhythmogram processes and output.*

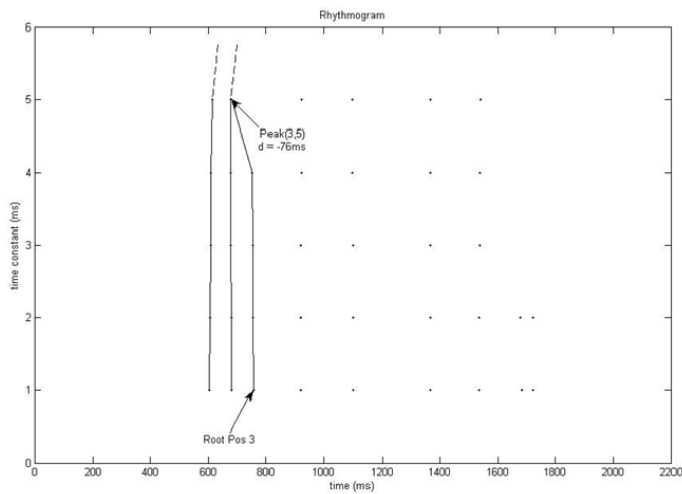


Figure 3.2: *Syllable-type event string within utterance.*

Table 3.1: Quantization of time distances between peaks.

Distance (ms)	Mapping	Distance (ms)	Mapping
0	1	-1	11
1	2	-2	12
2	3	-3	13
3	4	-4	14
4	5	-5	15
5 - 7	6	-6 - -8	16
8 - 10	7	-9 - -10	17
11 - 20	8	-11 - -20	18
21 - 50	9	-21 - -50	19
50 - ∞	10	-51 - -∞	20

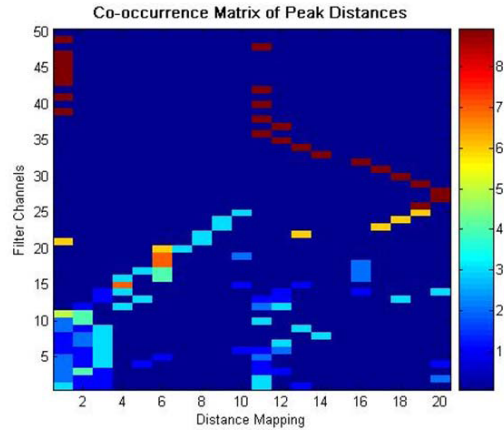


Figure 3.3: Co-occurrence matrix of quantized peak distances.

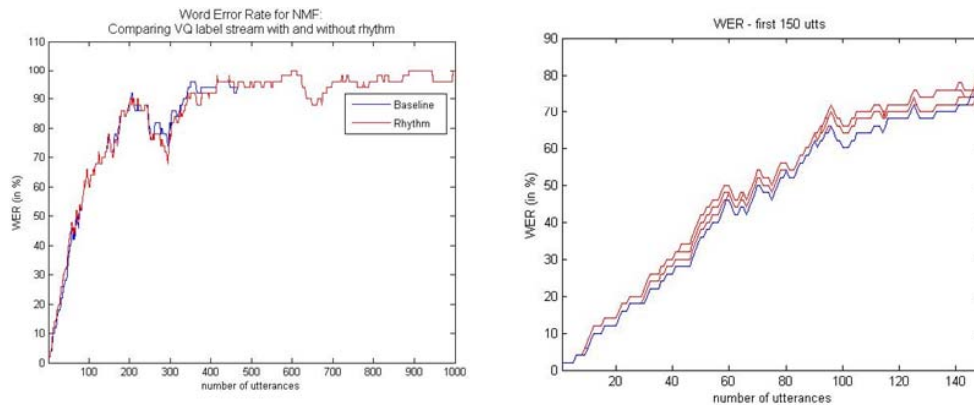


Figure 3.4: Plot of the key word tag recognition accuracy with and without rhythm a) for 1000 utterances b) for 150 utterances.

3.3 Rhythm VQ results

The first experiment carried out was to simply append the rhythm stream to the spectral input. Figure 3.4 shows the key word recognition accuracy using the same experimental set-up as described in section 3.4.3. The original result (baseline) is in blue and is an average of five attempts to correctly guess the associated key word tag of 1000 incoming utterances. The red plot is the average of five attempts with the rhythm stream appended to the current speech VQ labels. It appears that there is not much difference between the two, with most of the variance occurring in the first 500 utterances.

It was hypothesised that using the rhythm stream would help NMF learn key words faster. Looking at the first 150 utterances it can be seen that the prediction of key words is consistently better than the baseline for all five attempts (Fig. 3.4).

The next set of experiments was carried out with varying weights of the rhythm stream. The plots in Fig. 3.5 show the accuracy (%) difference from the baseline taken every 50 utterances for each weighting. The only weighting that is consistently better than the baseline is a factor of 0.1, where there is an improvement of nearly 2% after the first 100 utterances. The biggest accuracy difference

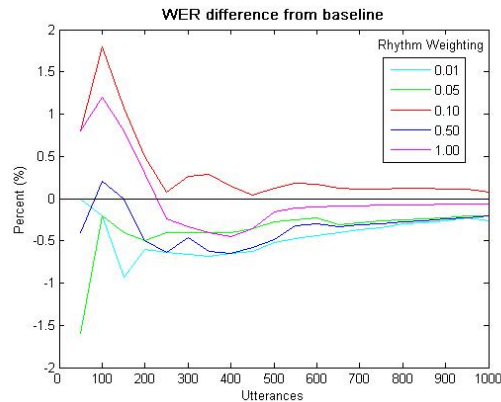


Figure 3.5: Accuracy difference from the baseline with different rhythm weightings for 1000 utterances.

is during the first 300 utterances, after this period they all tend towards baseline. Thus, the largest difference in accuracy is at the very beginning of the learning phase, but once the system has created a stable representation, then the rhythm stream does not improve performance.

3.4 Pitch Vectorisation

The second prosodic cue under investigation here is pitch. To vectorise the pitch for NMF we devised a procedure that relates pitch movement at different time instants within the utterance. This is achieved by calculating the deltas of the pitch contour and then accumulating the co-occurrence counts of a user defined lag value within this stream.

This procedure consists of three steps. In step 1, the pitch contour is calculated. There are two pitch extraction methods being used and compared in these experiments; the ACORNS *Pitch Estimator* (described in deliverable D1.1) and the *Subharmonic* function [2] which uses dynamic programming smoothing. Both methods carry out speech detection processes to give voiced/unvoiced values which are used to exclude unvoiced regions.

In step 2, the frequency range of the pitch contour within the voiced regions (i.e. regions consisting of consecutive voiced frames) is then quantized into a 50-channel filter bank. Fig. 3.6 shows the plot of the quantized pitch contour of the example utterance. The pitch contour was calculated either using the 'PitchEstimator' or the 'Subharmonic' method and then quantized to 50 frequency bins. Both methods calculate unvoiced regions of the signal which were removed.

We are now able to accumulate counts of lag- τ co-occurrences using:

$$C(q_t, q_{t-\tau}) = C(q_t, q_{t-\tau}) + 1, \quad (3.1)$$

where C is the co-occurrence matrix, q_t is the label of the quantised pitch stream at time t and τ is the lag of frames of 10 ms.

Figure 3.7 shows the plot of the co-occurrence matrix within the example utterance. This data matrix is appended to the NMF input stream. The figure shows that there is significant pitch variation as otherwise clusters of dots along the diagonal of the plot from top-left to bottom-right would be seen.

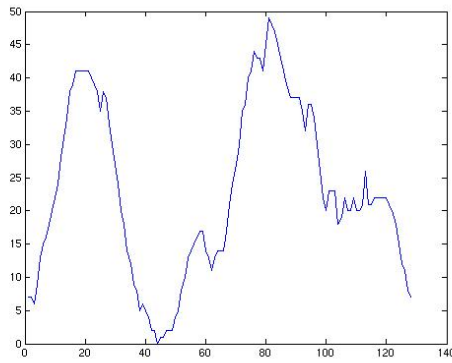


Figure 3.6: *Island of voiced pitch contour after quantization.*

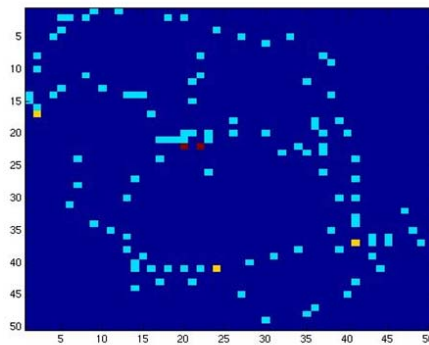


Figure 3.7: *A plot of the co-occurrence matrix showing pitch variation within the example utterance.*

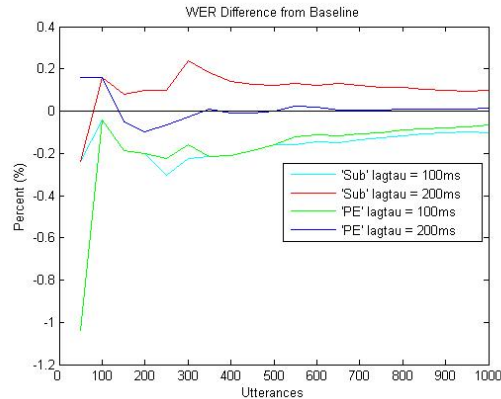


Figure 3.8: Accuracy difference from the baseline for pitch VQ using Subharmonic and pitchEstimator.

3.5 Pitch VQ results

The key word recognition accuracy difference from the baseline is plotted in Fig. 3.8.

The results show that using a longer lag-tau greater accuracy while only the *Subharmonic* method for calculating the pitch contour was better than the baseline.

3.6 Conclusions

As hypothesised, the addition of prosodic cues as an aid for word detection helped raise accuracy results during the early learning period. Rhythm had more of an impact than the use of the pitch contour with accuracy almost 2% better than the baseline during the first 100 utterances. After this period the accuracy tends towards the baseline for both cues. Calculating the pitch contour with dynamic programming smoothing slightly enhanced the results, but at the expense of computational complexity. Thus, the overall conclusion is that the improvements in word detection achieved by using rhythm are marginal.

Bibliography

- [1] N. P. M. Todd and G. J. Brown, "Visualization of rhythm, time and metre," *Artificial Intell. Rev.*, vol. 10, pp. 253–273, 1996.
- [2] T. J. Klasesen, "Robust pitch detection by subharmonic summation," pp. 253–273, 2003, available at: www.esat.kuleuven.be/psi/spraak/publications.

Chapter 4

Feature Selection Based on Knowledge of the Auditory System

C. Koniaris, M. Kuropatwinski, W.B. Kleijn (KTH)

4.1 Introduction

Recent improvements in performance of automatic speech recognition (ASR) systems can be attributed to a large extent to the development of effective acoustic modeling schemes. It is, however, generally accepted that the representation of the acoustic data is an important issue in the design and performance of any ASR system. In other words, if the speech features used for acoustic modeling do not include all relevant information available in the speech signal, then the performance of the classification stage is inherently suboptimal and likely cannot reach human recognition performance.

It is still an open issue whether all relevant information needed in distinguishing words is preserved by the front-end. In all cases, the goal is to reduce the dimensionality but often vital information of the original signal can be lost. The most widely used features in speech recognition, are the mel-frequency cepstral coefficients (MFCCs). The popularity of MFCCs among researchers is motivated by their low complexity and the high recognition rates especially for clean environments [1].

To reduce the dimensionality of given feature sets, algorithms have been proposed in the literature to select optimal subsets. One approach is to find the maximum statistical dependency between a feature subset and a class by computing the mutual information. This method is computationally intractable. An alternative approach proposed in [2], combines the minimal-redundancy-maximal-relevance (mRMR) criterion with a wrapper, a method to minimize the classification error for a particular classifier. In [3], the maximum entropy discrimination (MED) feature selection proposed for ASR. Results were comparable to a wrapper but the algorithm was less computationally expensive. In all methods, the relation between features and target classes was investigated and different criteria were applied to reduce the classification error.

The human hearing system has been modeled by complicated auditory models. Since the perceptually optimal processing of speech signal is difficult, a sedulous effort using distortion measures that are based on human perception is needed. In [4] the assumption of asymptotically small errors was used to construct more convenient distortion measures based on the so-called sensitivity matrix. This theme was later developed further in the context of rate-distortion theory [5, 6].

In [7], the sensitivity matrix is used to simplify a perceptual distortion measure for its use in audio

coding. In this paper, we use the sensitivity matrix to select features that humans perceive. In our approach we establish a measure of goodness for a given feature based on a perturbation analysis and distortion criteria derived from psycho-acoustic models. Based on this measure of goodness, a compact set of relevant features is derived. We assume that both the features and the distortion criteria are continuous and differentiable functions of the speech signal.

The paper is organized as follows. In Sec. 4.2 we introduce the distortion measure and the van de Par auditory model [8]. In Sec. 4.3, we present our algorithm for rating the perceptual significance of features. In Sec. 4.4 we apply the method in ASR using MFCCs as features. In Sec. 4.5 we investigate the range of the linearity assumption and show recognition results. Finally, in Sec. 4.6, we discuss our conclusions.

4.2 Preliminaries

We are interested in approximating quantitative models of human perception. We first discuss a distortion analysis and then its application to a particular auditory model.

4.2.1 Distortion analysis

In [4] the concept of the sensitivity matrix was introduced to approximate a given distortion measure used in the problem of quantization of the linear predictive coding (LPC) parameters in speech coding systems. Later, this work was extended and generalized in [5] and in [6]. In [7], a method for deriving the sensitivity matrix for distortion measures that are relevant for audio signals was developed based on spectro-temporal auditory models.

Let $\mathbf{x} = [x_1, x_2, \dots, x_N]^T$ be a speech signal vector. \mathbf{x} can be simply a time-domain vector, but it can also be, for example, a periodogram. Furthermore, let $\hat{\mathbf{x}}$ be a distorted version of \mathbf{x} and let $d[\mathbf{x}, \hat{\mathbf{x}}]$ be a distortion measure between \mathbf{x} and $\hat{\mathbf{x}}$. For small distortions, we perform a Taylor series expansion of d

$$d[\mathbf{x}, \hat{\mathbf{x}}] = d[\hat{\mathbf{x}}, \hat{\mathbf{x}}] + \left. \frac{\partial d[\mathbf{x}, \hat{\mathbf{x}}]}{\partial \hat{\mathbf{x}}} \right|_{\hat{\mathbf{x}}=\mathbf{x}} (\mathbf{x} - \hat{\mathbf{x}}) + \frac{1}{2} (\mathbf{x} - \hat{\mathbf{x}})^T \left. \frac{\partial^2 d[\mathbf{x}, \hat{\mathbf{x}}]}{\partial \hat{x}_i \partial \hat{x}_j} \right|_{\hat{\mathbf{x}}=\mathbf{x}} (\mathbf{x} - \hat{\mathbf{x}}) + \mathbf{O}[\|\mathbf{x} - \hat{\mathbf{x}}\|^3]. \quad (4.1)$$

In the above expansion we know that $d[\hat{\mathbf{x}}, \hat{\mathbf{x}}] = 0$, and because $\hat{\mathbf{x}}$ is a unique minimum of $d[\mathbf{x}, \hat{\mathbf{x}}]$, the term $\left. \frac{\partial d[\mathbf{x}, \hat{\mathbf{x}}]}{\partial \hat{\mathbf{x}}} \right|_{\hat{\mathbf{x}}=\mathbf{x}}$ vanishes. Moreover, all the terms that are of order three and above $\mathbf{O}[\|\mathbf{x} - \hat{\mathbf{x}}\|^3]$, are approximated to zero. Hence, the distortion measure is approximated [4] as

$$d[\mathbf{x}, \hat{\mathbf{x}}] \approx \frac{1}{2} (\mathbf{x} - \hat{\mathbf{x}})^T \mathbf{D}_{\mathbf{x}}[\mathbf{x}] (\mathbf{x} - \hat{\mathbf{x}}). \quad (4.2)$$

where

$$\mathbf{D}_{\mathbf{x}}[\mathbf{x}] = \left. \frac{\partial^2 d[\mathbf{x}, \hat{\mathbf{x}}]}{\partial \hat{x}_i \partial \hat{x}_j} \right|_{\hat{\mathbf{x}}=\mathbf{x}} \quad (4.3)$$

is called the *sensitivity matrix*. The word ‘‘sensitivity’’ refers to the fact that each element of this matrix represents the sensitivity of the distortion $d[\mathbf{x}, \hat{\mathbf{x}}]$ to a particular $(\mathbf{x} - \hat{\mathbf{x}})$.

4.2.2 van de Par auditory model

The van de Par [8] auditory model is a psychoacoustic masking model that accounts for simultaneous processing of sound signals. Let \mathbf{x} denote the square-root of the periodogram of a signal segment. The model consists of channels, which we index with f , in each of which the ratio of the distortion $\mathbf{x}(f) - \hat{\mathbf{x}}(f)$ to masker $\mathbf{x}(f)$ is estimated. In the end, all ratios are combined together, to account for the spectral integration property of the human auditory system. The complete model is then described by

$$d[\mathbf{x}, \hat{\mathbf{x}}] = C_s L_{\text{eff}} \sum_i \frac{\frac{1}{N} \sum_f |\mathbf{h}_{\text{om}}(f)|^2 |\gamma_i(f)|^2 |\mathbf{x}(f) - \hat{\mathbf{x}}(f)|^2}{\frac{1}{N} \sum_f |\mathbf{h}_{\text{om}}(f)|^2 |\gamma_i(f)|^2 |\mathbf{x}(f)|^2 + C_a}, \quad (4.4)$$

where C_s and C_a are constants calibrated based on measurement data, L_{eff} is the effective duration of the segment according to the temporal integration time of the human auditory system, $\mathbf{h}_{\text{om}}(f)$ is the outer and middle ear transfer function which is the inverse of the threshold in quiet and finally $\gamma_i(f)$ is the i 'th gammatone filter.

In our system, the van de Par model is used to obtain the sensitivity matrix in the speech frequency domain. Combining (4.3) and (4.4), we obtain that $\mathbf{D}_{\mathbf{x}}[\mathbf{x}]$ is a diagonal matrix with diagonal element

$$\mathbf{D}_{\mathbf{x},ff}[\mathbf{x}] \approx 2C_s L_{\text{eff}} \sum_i \frac{\frac{1}{N} |\mathbf{h}_{\text{om}}(f)|^2 |\gamma_i(f)|^2}{\frac{1}{N} \sum_f |\mathbf{h}_{\text{om}}(f)|^2 |\gamma_i(f)|^2 |\mathbf{x}(f)|^2 + C_a} \quad (4.5)$$

for row f .

4.3 A method for rating the perceptual significance of features

The relative importance of perturbation vector and, therefore, the corresponding feature to a psychoacoustic distortion measure can be established with a sensitivity analysis of the distortion measure for a given speech vector. A sensitivity analysis [7] establishes a basis for the speech block where the basis vectors have ordered sensitivity. That is, for a given speech vector we know which perturbation vector in the signal space has the largest impact on distortion, which perturbation vector in the remaining subspace (signal space without the first perturbation vector) has the largest impact on distortion, etc. Thus, by setting an audibility threshold, we can establish a signal subspace where changes to the signal are most audible. We call this subspace the perceptually relevant subspace. Combining the above feature and distortion measure analysis, we derive compact signal representations based on human perception.

4.3.1 Feature and speech: a linearized relation

The relationship between the extracted features and the speech signal is not linear in general. However, we can linearize it around the observed speech vector \mathbf{x} . We use as speech vector the square root of the periodogram, since both the features we study and the van de Par model are a function of

these magnitude spectra. Let $\mathbf{c}[\mathbf{x}]$ be the transform of the spectral description of the speech segment onto the feature vector i.e., $\mathbf{c} : \mathbb{R}^N \rightarrow \mathbb{R}^Q$. Then we can write

$$\mathbf{c}[\hat{\mathbf{x}}] \approx \mathbf{c}[\mathbf{x}] + \left. \frac{\partial \mathbf{c}[\hat{\mathbf{x}}]}{\partial \hat{x}} \right|_{\hat{\mathbf{x}}=\mathbf{x}} (\hat{\mathbf{x}} - \mathbf{x}). \quad (4.6)$$

Denoting $\mathbf{A} = \left. \frac{\partial \mathbf{c}[\hat{\mathbf{x}}]}{\partial \hat{x}} \right|_{\hat{\mathbf{x}}=\mathbf{x}}$, where $\mathbf{A} \in \mathbb{R}^{Q \times N}$, we arrive at

$$\delta \mathbf{c} \approx \mathbf{A} \delta \mathbf{x}, \quad (4.7)$$

where $\delta \mathbf{c} = \mathbf{c}[\mathbf{x}] - \mathbf{c}[\hat{\mathbf{x}}]$ and $\delta \mathbf{x} = \mathbf{x} - \hat{\mathbf{x}}$.

Eq. (4.7) gives a linear approximation of the transformation of the speech error vector to the feature. Since the above equation is underdetermined, it has infinite set of solutions. However, all solutions if projected onto the image subspace of \mathbf{A} give the same result. This projected solution results in a N -dimensional vector $\delta \mathbf{x}_c$ lying in Q -dimensional subspace and has the following form

$$\delta \mathbf{x}_c \approx \mathbf{A}^+ \delta \mathbf{c}. \quad (4.8)$$

In the next section, we show how to use this result to compute the sensitivity matrix in the feature space.

4.3.2 The sensitivity matrix in the feature domain

In Sec. 4.2.1, we derived the distortion measure shown that it can be expressed as in (4.3). Moreover, in Sec. 4.2.2, we obtained the sensitivity matrix (4.5) in the speech frequency domain from the van de Par auditory model.

We can now move to the feature domain, and use (4.8) to compute the new sensitivity matrix for the features

$$d[\mathbf{c}, \hat{\mathbf{c}}; \mathbf{x}] = \frac{1}{2} \delta \mathbf{x}_c^T \mathbf{D}_x \delta \mathbf{x}_c = \frac{1}{2} \delta \mathbf{c}^T (\mathbf{A}^+)^T \mathbf{D}_x \mathbf{A}^+ \delta \mathbf{c} \quad (4.9)$$

or,

$$d[\mathbf{c}, \hat{\mathbf{c}}; \mathbf{x}] = \frac{1}{2} \delta \mathbf{c}^T \mathbf{D}_c[\mathbf{x}] \delta \mathbf{c}. \quad (4.10)$$

Matrix $\mathbf{D}_c[\mathbf{x}]$ is the new sensitivity matrix in the feature domain and can be considered that it describes a linearization of the perceptual transform [7].

Motivated by the assumption that good features subset should cover the perceptually relevant signal subspace, we define a measure of features goodness $G(i)$ to rate the perceptual significance of features set i over all speech segments j . We want the squared error in the feature domain to be proportional to the auditory-model distortion, for small distortion, and over all signal segments. A suitable measure is

$$G(i) = \min_i \left\{ \sum_j \left(\Gamma_j - \frac{\sum_k \Gamma_k \mathbf{B}_k(i)}{\sum_k \mathbf{B}_k(i) \mathbf{B}_k(i)} \mathbf{B}_j(i) \right)^2 \right\} \quad (4.11)$$

where we have introduced Γ_j as

$$\Gamma_j = \int_{\varepsilon} \delta \mathbf{x}^T \mathbf{D}_x \delta \mathbf{x} \, d\mathbf{x} \quad (4.12)$$

and $\mathbf{B}_j(i)$ as

$$\mathbf{B}_j(i) = \int_{\varepsilon} \delta \mathbf{x}^T \mathbf{A}^T(i) \mathbf{A}(i) \delta \mathbf{x} \, d\mathbf{x}. \quad (4.13)$$

In (4.12) and (4.13), ε is a small positive real number defining a region in which we compute distortions over a set of equal but very small norm errors for all speech signal segments. (In an implementation the value of ε can be set to the smallest value that facilitates reasonable computational precision.) During the computation of $G(i)$ both $\mathbf{A}(i)$ and \mathbf{D}_x are adapted to the specific segment \mathbf{x} . A small $G(i)$ means that the local distortion (averaged over small deviations) in the feature domain, \mathbf{B}_j , behaves like the local distortion (averaged over small deviations) in the auditory domain Γ_j except for a scaling, and hence the best features set should minimize the shape difference in (4.11).

4.4 Application to speech recognition

As we have already mentioned in the introduction, MFCCs are the features that most researchers use for speech recognition. Mel frequencies are based on the knowledge of the human auditory system. Human ear resolves frequencies in a nonlinear manner. The response is linear at frequencies below 1 kHz and becomes logarithmic with increasing frequency [1].

4.4.1 The \mathbf{A} matrix for MFCCs

Mel-frequency cepstrum coefficients (MFCCs) can be computed as [1]

$$\mathbf{c}[q] = \sum_{m=0}^{M-1} \mathbf{s}[m] \cos \left[q \left(m - \frac{1}{2} \right) \frac{\pi}{M} \right], \quad q = 1, 2, \dots, Q, \quad (4.14)$$

where Q is the number of cepstrum coefficients, $\mathbf{s}[m]$ represents the log-energy output of the m 'th filter of the filterbank, and M denotes the number of triangular bandpass filters used.

In Sec. (4.3.1) we introduced the matrix \mathbf{A} that characterizes the local relation between the features and the signal \mathbf{x} . To find \mathbf{A} for MFCCs, we follow the steps on computing them, backwards. As a result, (4.14) can then be written as

$$\mathbf{c}[q] = \sum_{m=0}^{M-1} \ln \mathbf{z}[m] \cos \left[q \left(m - \frac{1}{2} \right) \frac{\pi}{M} \right], \quad (4.15)$$

where $\mathbf{z}[m]$ is the product of power spectrum and the triangular mel weighted filters or,

$$\mathbf{c}[q] = \sum_{m=0}^{M-1} \ln \left\{ \sum_{n=0}^{N-1} \mathbf{x}[n] \mathbf{H}_m[n] \right\} \cos \left[q \left(m - \frac{1}{2} \right) \frac{\pi}{M} \right], \quad (4.16)$$

where $\mathbf{x}[n]$ is the periodogram and $\mathbf{H}_m[n]$ is the m 'th triangular mel-filter. From the above, we can calculate \mathbf{A} as the product of the following derivatives

$$\mathbf{A}[q, n] = \frac{\partial \mathbf{c}[q]}{\partial \mathbf{s}[m]} \frac{\partial \mathbf{s}[m]}{\partial \mathbf{z}[m]} \frac{\partial \mathbf{z}[m]}{\partial \mathbf{x}[n]}, \quad (4.17)$$

which is

$$\mathbf{A}[q, n] = \sum_{m=0}^{M-1} \cos \left[q \left(m - \frac{1}{2} \right) \frac{\pi}{M} \right] \frac{1}{\mathbf{z}[m]} \mathbf{H}_m[n]. \quad (4.18)$$

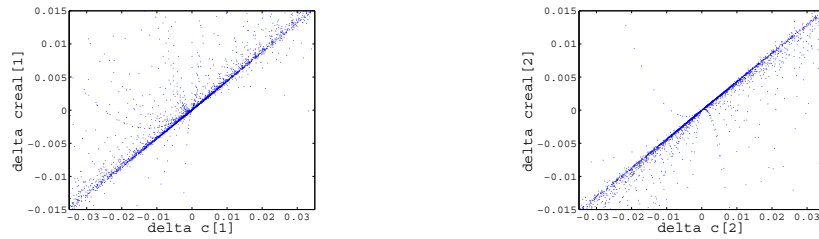


Figure 4.1: Scatter plots of the estimated δc 's vs. δc_{true} 's for the 1st and 2nd MFCC, respectively.

4.4.2 The selection process

The goal is to find the subset of MFCCs that describes the most audible signal components best, in an average sense. To this purpose we introduced the *goodness gain* G in the previous section. The process starts by computing in every segment, the value of Γ_j and the sensitivity matrix \mathbf{D}_x . Then, for each possible combination of features subsets i , we consider the matrix $\mathbf{A}(i)$ and calculate the quantity $\mathbf{B}_j(i)$. We repeat the above procedure for all the speech segments. In the end, we compute the goodness gain measure for every feature subset and for the whole speech. A comparison of the goodness gain values of all possible subsets is considered, and the one that has the minimum value is rated to be the optimal in terms of minimizing the shape difference as it is given by (4.11).

4.5 Evaluation

4.5.1 Range of linearization

In this we examine the range of the linearization assumption between the cepstrum and the speech. The speech is distorted with i.i.d. Gaussian noise at different SNRs ranging from 30 to 90 dB. The sampling frequency is at 16 KHz, and the segments length is 25 ms.

Test Set A	12 MFCC, C0,E		11 MFCC, C0,E		10 MFCC, C0,E		9 MFCC, C0,E		8 MFCC, C0,E	
	full set	AMFS	ref	AMFS	ref	AMFS	ref	AMFS	ref	
Clean 1	97.6 %	97.6 %	97.3 %	97.9 %	97.1 %	97.7 %	95.5 %	97.6 %	94.8 %	
Clean 2	97.1 %	97.0 %	96.8 %	97.3 %	96.5 %	96.9 %	95.0 %	97.0 %	94.5 %	
Clean 3	97.3 %	97.3 %	96.9 %	96.9 %	96.6 %	97.0 %	95.1 %	96.9 %	94.6 %	
Clean 4	97.5 %	97.7 %	97.1 %	97.7 %	96.7 %	97.7 %	95.2 %	97.4 %	94.8 %	
Average	97.4 %	97.4 %	97.0 %	97.5 %	96.7 %	97.3 %	95.2 %	97.2 %	94.7 %	

Test Set A	7 MFCC, C0,E		6 MFCC, C0,E		5 MFCC, C0,E		4 MFCC, C0,E		3 MFCC, C0,E	
	AMFS	ref	AMFS	ref	AMFS	ref	AMFS	ref	AMFS	ref
Clean 1	97.4 %	93.4 %	97.3 %	90.4 %	97.1 %	88.8 %	96.9 %	84.0 %	86.5 %	79.8 %
Clean 2	97.0 %	93.0 %	96.9 %	90.2 %	96.4 %	88.6 %	96.3 %	84.3 %	86.2 %	80.0 %
Clean 3	96.9 %	93.4 %	96.6 %	90.3 %	96.7 %	88.7 %	96.5 %	84.4 %	86.6 %	79.8 %
Clean 4	97.5 %	93.2 %	97.1 %	90.6 %	97.0 %	88.6 %	96.8 %	84.2 %	86.6 %	79.4 %
Average	97.2 %	93.3 %	97.0 %	90.4 %	96.8 %	88.7 %	96.6 %	84.2 %	86.5 %	79.8 %

Table 4.1: Recognition accuracy for the full set of 12 MFCCs and for several subsets of auditory-model based feature selection (AMFS). The reference (ref) is the average accuracy obtained from 5 different, randomly selected features subsets.

Fig. 4.1 shows the δc i.e., computed from the linearized relation (4.7) versus the true difference δc_{true} between the cepstrum of the original signal and the cepstrum of the distorted one for only the

1st and 2nd MFCC, respectively. It can be seen, that the linearity assumption is a good approximation in this range of SNRs.

4.5.2 Speech recognition experiments

In this section we present experimental results comparing the selected feature set, obtained from our proposed method, to the conventional feature sets that are used in most cases.

We use the AURORA 2 [9] database. The MFCCs are extracted by using a Hamming window of 25 ms with an overlap of 12.5 ms. The length of the DFT is set to 256, while the number of filters used are 23. In the end, a set of 12 conventional MFCCs are extracted. As a recognizer we use the HTK [10] toolbox. The digits are modeled as whole word HMMs with 16 states (HTK notation is 18 states including the beginning and end states) and three Gaussian mixture components per state. At first, we produce an initial model with global data means and variances, same for each digit and then we run 16 iterations to build the final model. Table 4.1 shows the recognition accuracy for only the

subset	MFCCs
11	c[1],c[2],c[3],c[4],c[5],c[7],c[8],c[9],c[10],c[11],c[12]
10	c[1],c[2],c[3],c[4],c[5],c[8],c[9],c[10],c[11],c[12]
9	c[1],c[2],c[3],c[4],c[5],c[9],c[10],c[11],c[12]
8	c[1],c[2],c[3],c[4],c[9],c[10],c[11],c[12]
7	c[1],c[2],c[3],c[4],c[10],c[11],c[12]
6	c[1],c[2],c[3],c[4],c[11],c[12]
5	c[1],c[2],c[3],c[11],c[12]
4	c[1],c[2],c[3],c[12]
3	c[1],c[11],c[12]

Table 4.2: Selected MFCCs

clean part of test set A of the database. In the first column, the recognition accuracy of the full set is shown and in the next columns, the accuracy of the auditory-mobel based feature selection (AMFS) in different cardinalities. For comparison, the average performance of 5 different, randomly selected MFCCs subsets is shown, too. The performance of the AMFS remains sufficiently stable, comparable to the full set, as the number of coefficients reduces. In comparison to the reference configuration, AMFS has better recognition accuracy in a range from 0.4% up to 12.4%. Finally, table 4.2 shows the selected MFCCs. The system favours not to choose coefficients from the middle part of the MFCCs' range, keeping only the initial and last coefficients.

4.6 Conclusions

We presented a new method to select speech features based on human perception. The selection algorithm was based on a particular, relatively simple auditory model. We applied it to MFCCs to find an optimal set that removes redundancy and thus, lowers dimensionality. We evaluated the subsets with a series of classification experiments on the AURORA2 speech database. Results showed that the system can indeed perform well based only on perception, and ignores the aspects of sound that we do not hear. In the future, it is natural to apply the method to auditory models that include the effect of time-domain masking.

Bibliography

- [1] S. B. Davis and P. Mermelstein, "Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences," *IEEE Trans. Acoust., Speech, Sig. Proc.*, vol. 28, no. 4, pp. 357–366, Aug. 1980.
- [2] H. Peng, F. Long, and C. Ding, "Feature selection based on mutual information: Criteria of max-dependency, max-relevance, and min-redundancy," *IEEE Trans. Pat. Analys., Mach. Intellig.*, vol. 27, no. 8, pp. 1226–1238, Aug. 2005.
- [3] F. Valente and C. Wellekens, "Maximum entropy discrimination (MED) feature subset selection for speech recognition," *IEEE Workshop on ASRU*, pp. 327–332, Dec. 2003.
- [4] W. R. Gardner and B. D. Rao, "Theoretical analysis of the high-rate vector quantization of LPC parameters," *IEEE Trans. Speech, Audio Proc.*, vol. 3, no. 5, pp. 367–381, Sep. 1995.
- [5] T. Linder, R. Zamir, and K. Zeger, "High-resolution source coding for non-difference distortion measures: multidimensional companding," *IEEE Trans. Inform. Theory*, vol. 45, no. 2, pp. 548–561, Mar. 1999.
- [6] J. Li, N. Chaddha, and R. M. Gray, "Asymptotic performance of vector quantizers with a perceptual distortion measure," *IEEE Trans. Inform. Theory*, vol. 45, no. 4, pp. 1082–1091, May 1999.
- [7] J. H. Plasberg and W. B. Kleijn, "The sensitivity matrix: Using advanced auditory models in speech and audio processing," *IEEE Trans. Speech, Audio Proc.*, vol. 15, no. 1, pp. 310–319, Jan. 2007.
- [8] S. van de Par, G. Charestan, and R. Heusdens, "A gammatone-based psychoacoustical modeling approach for speech and audio coding," *Proc.ProRISC,Veldhoven,NL*, pp. 321–326, 2001.
- [9] H. G. Hirsch and D. Pearce, "The AURORA experimental framework for the performance evaluation of speech recognition systems under noise conditions," in *ISCA ITRW ASR2000*, Paris, 2000.
- [10] S. Young, D. Kershaw, J. Odell, D. Ollason, V. Valtchev, and P. Woodland, *The HTK Book (for HTK Version 3.2)*. Cambridge University, Engineering Department, Dec. 2002.