

# Noise Robust Feature Extraction for ASR using the Aurora 2 Database

Qifeng Zhu, Markus Iseli, Xiaodong Cui, Abeer Alwan

Department of Electrical Engineering  
University of California, Los Angeles  
{qifeng, iseli, xdcui, alwan}@icsl.ucla.edu

## Abstract

Four front-end processing techniques developed for noise robust speech recognition are tested with the Aurora 2 database. These techniques include three previously published algorithms: variable frame rate analysis [Zhu and Alwan, 2000], peak isolation [Strope and Alwan, 1997], and harmonic demodulation [Zhu and Alwan, 2000], and a new technique for peak-to-valley ratio locking. Our previous work has focused on isolated digit recognition. In this paper, these algorithms are modified for recognition of connected digits. Recognition results with the Aurora 2 database show that a combination of these four techniques results in 53% and 12% error rate reduction for the clean training and multicondition training, respectively, when compared to the baseline MFCC front-end, with no significant increase in computational complexity.

## 1. Introduction

This paper focuses on front-end feature extraction approaches for noise robust automatic speech recognition (ASR). Four front-end processing techniques are tested with the Aurora 2 database. These techniques include three previously published algorithms: variable frame rate analysis [3], peak isolation [2], and harmonic demodulation [4], and a new technique for peak-to-valley ratio locking. Our previous work has focused on isolated digit recognition and mainly computer-generated additive noise.

Here, training and testing followed the specifications described in [1]. A word-based ASR system for digit string recognition where each HMM word model has 16 emitting states is adopted. Training is done with either 8440 clean utterances (referred to as clean training) or with 8440 clean and noisy utterances (multi-condition training). A 3-state silence model and a one state short pause model are used. Test data included different kinds of realistic background noise at various SNRs.

The Aurora 2 database CD included a program (FE2.0) to compute the MFCCs and log energy. The front-end used by HTK is MFCC\_E\_D\_A, which contains 12 MFCCs and log energy together with their first and second derivatives. Each feature vector thus contains 39 components. Reference recognition results are computed with FE2.0. The techniques used in this paper were implemented by modifying the code in FE2.0.

## 2. Noise robust front-end feature extraction

Three previously published front-ends and a new algorithm are described in this section.

### 2.1. Variable frame rate analysis (VFR) [3]

Variable frame rate (VFR) analysis in [3] is motivated by the fact that changes in spectral characteristics are important cues for discriminating and identifying speech sounds. These changes can occur over very short time intervals. Computing frames every 10 ms, as commonly done in ASR, is not sufficient to capture such dynamic changes. The VFR algorithm increases the frame rate for rapidly-changing segments with relatively high energy and decreases the frame rate for steady-state segments, based on a weighted log energy Euclidean MFCC distance. The smallest frame step can be 2.5 ms. An example is shown in Figure 1. The current implementation uses an average frame rate which is less than 100 frames per second.

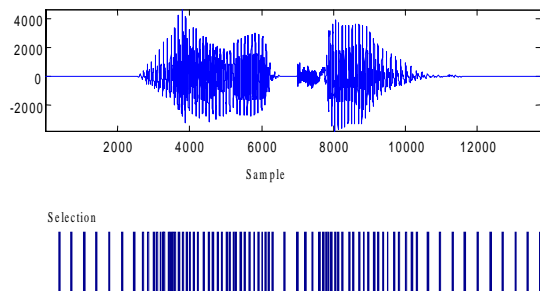


Figure 1: The upper panel shows the utterance "one two". The lower panel shows the selected frames.

MFCCs computed with the VFR technique reduce the error rate, when compared to reference results, in the clean training condition by 27.69% for Set A, 36.99% for Set B, and 0.92% for Set C. The overall error reduction is 27.27%.

### 2.2. Peak isolation (PKISO) [2]

This technique involves MFCC liftering, an inverse DCT (IDCT), and half wave rectification. After the IDCT, spectral valleys are often less than 0 and formants are larger than 0. Figure 2 shows an example of the recovered log Mel filter-bank output from liftered MFCCs for a clean and noisy frame of /i/. Half-wave rectification is then applied on the recovered log Mel filter-bank output so that the valleys are effectively removed. A DCT is applied on the rectified log Mel filter-bank output to obtain feature vectors which will be referred to as PKISO\_MFCCs.

Error rate reductions with PKISO\_MFCCs in the clean training condition are, when compared to reference results, 45.77% for Set A, 49.54% for Set B, and 10.60% for Set C. The overall error reduction is 41.48%.

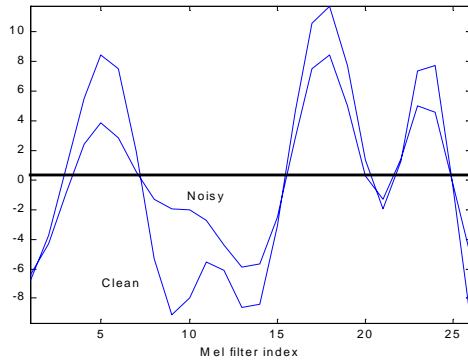


Figure 2: Peak isolation. Log Mel filter-bank output recovered from liftered MFCCs for a clean and noisy (0 dB SNR) frame of /i/. After half-wave rectification only the upper part is retained.

### 2.3. Peak-to-valley ratio locking

We introduce in this paper the concept of peak-to-valley ratio locking. In the presence of noise, spectral valleys will be buried by noise, but formants are, on average, not affected as much. An example can be seen in Figure 3, which shows spectra of a clean and noisy speech frame. The frame is from /i/ in “zero” (female talker) and the additive noise is speech shaped at 0 dB SNR. The noisy spectrum is an average over 150 frames. Note that the spectra are nearly the same at harmonic peaks around the formants, where the amplitude is about 3 times higher than the average noise spectrum. At frequencies where the signal amplitude is low, as in spectral valleys, the average noisy speech spectrum is nearly the same as the average noise spectrum. Because of the difference of the noise effects on valleys and peaks, often the peak-to-valley ratio in spectra of noisy speech is lower than that in clean speech, hence leading to a mismatch between the clean and noisy data.

After obtaining the recovered log Mel filter output from liftered MFCCs (without  $C_0$ ), as shown in Figure 2, both peaks and valleys will be affected. One approach to addressing this problem is peak-to-valley ratio locking. We set the highest peak of amplitude  $x$  to a fixed number  $p$ . The entire recovered Mel filter output is then scaled accordingly by a factor of  $p/x$ . In our implementation  $p$  was set to 10. This number is approximately the average amplitude of the highest peaks across the database.

When used together with peak isolation, only the positive part in the recovered Mel filter output is scaled, and the negative part is set to zero. An example of the result of combining PKISO with peak-to-valley ratio locking is shown in Figure 4.

The error rate reduction with peak-to-valley ratio locking only (without PKISO) in the clean training condition are: 37.32% for Set A, 42.22% for Set B, and 3.25% for Set C. The overall error reduction is 32.61%.

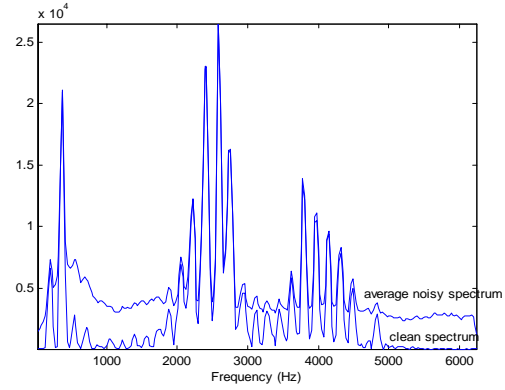


Figure 3: Linear clean and noisy (0 dB SNR) speech spectra.

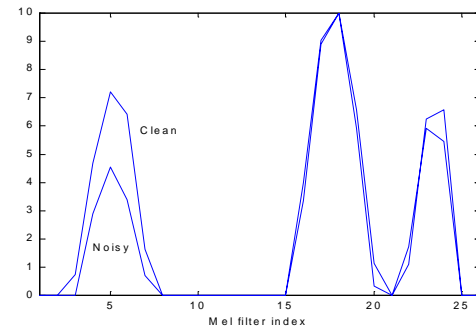


Figure 4: Log Mel filter-bank output after liftering, rectification and peak-to-valley ratio locking for the same frame of /i/ in Figure 2. Notice that the highest peaks are set to 10.

### 2.4. Harmonic demodulation (HD) [4]

Harmonic demodulation is a method that aims at reducing the difference between clean and noisy speech spectra especially at inter-harmonic valleys.

The LTI speech production model is viewed as amplitude modulation in the frequency domain, with the excitation spectrum being the carrier and the spectrum of the vocal tract transfer function being the modulator. Non-coherent demodulation with non-linear envelope detection is used to recover the spectrum of the vocal tract transfer function [4]. Envelopes of the speech spectra, instead of the speech spectra themselves, are used to compute the MFCCs.

The error rate reduction with harmonic demodulation in the clean training condition are: 28.37% for Set A, 33.21% for Set B, and 5.62% for Set C. The overall error reduction is 26.66%.

### 3. Modifications of the algorithms for the Aurora 2 database

#### 3.1. Speech/nonspeech detection

Our previous studies on VFR [3] PKISO [2], and HD [4], focused on isolated digits with endpoint detection. All techniques mentioned in this paper assume a speech model, hence speech/nonspeech detection is critical.

Speech/nonspeech detection is based on energy and voicing. The general idea is that if a frame has high energy and is a voiced frame or is close to voiced frames, then it is classified as speech. The energy condition is implemented by comparing the log-energy of a frame to a threshold for the corresponding utterance. A threshold (T) for an utterance is determined empirically by  $T=(H+L)*0.5$ , where H and L are the average of the 10 highest and 10 lowest log energy values in the utterance, respectively. The voicing detection is performed with the open source software from the SVR group at University of Cambridge, Pitch\_tracker 1.0, which is based on [6]. If the log energy of a frame is higher than a threshold and it is voiced or it is within 30 ms from a voiced frame, then this frame is classified as speech and the algorithms are applied.

The VFR algorithm is applied on each speech segment longer than 30 ms with a minimum frame step of 2.5 ms. For any speech segment shorter than 30 ms, a 10 ms frame step is used. A frame step of 25 ms is used for nonspeech segments.

#### 3.2. Increasing the variances of the silence model

The HD, PKISO and peak-to-valley ratio locking algorithms remove the mean spectral difference between the clean and noisy speech spectra. For the silence model, however, which is trained with clean data, there will be a large mismatch with the noisy test data. One solution is to increase the variances in the silence and short pause models.

We found that if the reference MFCCs are used, for the clean training condition, best performance is achieved by increasing the variances of the silence model by a factor of 1.1, the overall improvement in accuracy rate is less than 2%. With PKISO, HD, and peak-to-valley ratio locking, we get a better model of the digits, and hence a silence model with larger variance results in a larger increase in recognition performance. The increasing factor we use is 1.2, and the improvement in overall error rate reduction is 4.6%.

#### 3.3. Rasta-like filter in the cepstral domain

The four techniques described in this paper have difficulty with Set C, where the problem is channel mismatch. A Rasta like band pass filter [5] is used at the final stage of the front-end to avoid the mean shift effect caused by channel distortion. This results in a 5% improvement in error rate reduction in Set C, and a 0.5% and 1.3% improvement for sets A and B, respectively, compared with the results without the filter.

### 4. Complexity considerations

Even though our implementation did not optimize for computational cost (but focused on optimizing recognition performance) the increase in computational complexity of PKISO, HD and peak-to-valley ratio locking together is not high. In addition, these algorithms are frame-based and thus do not introduce a delay. The speech/non-speech detection algorithm introduces a delay equal to the utterance duration.

The additional memory cost for HD is 128 floating numbers, which is half of the FFT length and the additional memory for PKISO and peak-to-valley ratio locking together is 23 floating numbers, which is the number of the Mel filters.

The extra operations introduced for processing one frame in HD mainly comes from  $7*128$  (7 is the length of the filter characteristics and 128 is half the FFT size) floating number multiplications. The extra operations introduced by PKISO and peak-to-valley ratio locking mainly come from the extra IDCT and DCT, which contain  $23*12$  (23 is the number of the Mel filters and 12 is the length of the MFCC vector) floating number multiplication each. Liftering in PKISO adds 12 multiplications and peak-to-valley ratio locking adds another 23 multiplications in processing one frame.

When tested on a Sun Ultra Sparc 60 workstation, the computation load of PKISO and peak-to-valley ratio locking together is less than 4% of the total computation time of the original front-end (FE2.0) executable. HD adds about 20% more computational time. These measurements only count the front-end computation time, excluding the disk I/O time.

The computational load of VFR is higher than the other algorithms and depends on the number of frames classified as speech in an utterance. A delay equivalent to the duration of speech segments is introduced. The threshold for frame selection is computed from the inter-frame MFCC distance.

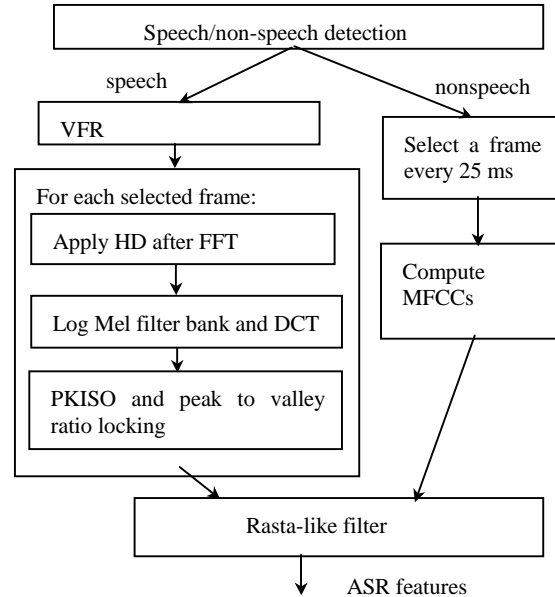


Figure 5: Block diagram of the front-end processing for clean training.

Aurora 2 Clean Training - Results														Percentage Improvement	
	A					B				C			Overall		
	Subwav	Babble	Car	Exhibition	Average	Restaurant	Street	Airport	Station	Average	Subway M	Street M		Average	
Clean	98.62	98.31	98.42	98.49	98.46	98.62	98.31	98.42	98.49	98.46	98.28	98.04	98.16	98.40	-66.76%
20 dB	96.68	97.25	97.20	96.76	96.97	97.18	96.95	97.32	97.50	97.24	94.87	95.74	95.31	96.75	32.00%
15 dB	94.50	96.22	95.65	94.60	95.24	95.64	95.71	95.85	95.77	95.74	91.19	93.14	92.17	94.83	56.17%
10 dB	89.62	92.90	91.38	87.87	90.44	91.25	91.14	91.02	90.65	91.02	82.04	85.19	83.62	89.31	64.52%
5 dB	76.33	83.01	81.93	75.10	79.09	75.47	79.81	80.55	77.23	78.27	59.66	69.26	64.46	75.84	58.36%
0 dB	47.74	53.90	61.50	47.58	52.68	44.67	55.86	58.51	56.56	53.90	28.49	39.69	34.09	49.45	38.20%
-5dB	19.56	19.01	23.65	17.00	19.81	12.25	22.94	22.49	22.93	20.15	10.29	14.90	12.60	18.50	10.80%
Average	80.97	84.66	85.53	80.38	82.89	80.84	83.89	84.65	83.54	83.23	71.25	76.60	73.93	81.23	
	37.65%	69.38%	63.28%	43.31%	55.73%	59.59%	58.15%	67.17%	62.91%	62.11%	15.03%	30.95%	23.00%		53.01%

Table 4.1: Clean training results on Aurora 2 database using all four techniques mentioned in this paper.

Aurora 2 Multicondition Training - Results														Percentage Improvement	
	A					B				C			Overall		
	Subwav	Babble	Car	Exhibition	Average	Restaurant	Street	Airport	Station	Average	Subway M	Street M		Average	
Clean	98.62	98.70	98.75	98.98	98.76	98.62	98.70	98.75	98.98	98.76	98.56	98.55	98.56	98.72	12.67%
20 dB	98.22	98.04	98.09	98.24	98.15	98.22	98.25	98.18	98.61	98.32	97.30	97.58	97.44	98.07	25.74%
15 dB	96.75	97.19	97.67	97.19	97.20	96.78	96.98	97.49	97.35	97.15	95.61	96.28	95.95	96.93	15.33%
10 dB	94.32	95.28	94.96	93.83	94.60	94.63	94.77	94.99	95.68	95.02	91.46	92.47	91.97	94.24	5.15%
5 dB	89.68	89.69	88.25	84.97	88.15	87.66	88.09	88.37	87.44	87.89	79.80	80.96	80.38	86.49	6.87%
0 dB	73.26	69.17	64.57	65.01	68.00	66.50	68.95	72.11	65.69	68.31	46.91	50.73	48.82	64.29	13.48%
-5dB	31.62	33.98	22.28	24.31	28.05	31.99	30.65	35.52	25.70	30.97	19.25	22.07	20.66	27.74	4.39%
Average	90.45	89.87	88.71	87.85	89.22	88.76	89.41	90.23	88.95	89.34	82.22	83.60	82.91	88.00	
	15.03%	15.97%	16.21%	-1.52%	11.52%	23.04%	18.28%	20.93%	26.31%	22.34%	-6.14%	-4.50%	-5.34%		11.86%

Table 2: Multicondition training results on Aurora 2 database, only harmonic demodulated and VFR are used in computing MFCCs.

## 5. Recognition results

Recognition experiments were performed with scripts included in the Aurora 2 CD, and with HTK 2.2. Training follows the steps specified in [1]. The four techniques were combined to produce the front-end for ASR as shown in Figure 5. Tables 1 and 2 show the results (word accuracy) for the clean and multi-condition training, respectively, when tested with sets A, B, and C at seven SNRs. Improvements in error reduction, when compared to the reference results with MFCCs, are shown in the last row and rightmost column in each table.

As Table 1 shows, the methods mentioned in this paper reduce the overall error rate by 53.01%. Error reductions are 55.7%, 62.1% and 23% for Sets A, B, and C, respectively.

VFR has high computational complexity. If VFR is not used in combination with the other algorithms, the overall error reduction is 48.82%. Error reductions are 52.8% for Set A, 55.07% for Set B, and 23.4% for Set C.

As shown in Table 1, the algorithms degrade slightly the performance of the clean condition (from 99% to 98.4%). But for all other SNRs the performance is improved significantly. Error reduction is best with babble (70%), airport (68%), car (63%), and train station (63%) background noise.

Most of the techniques reported in this paper (PKISO, peak-to-valley ratio locking, and HD) aim at reducing the difference of the means between clean and noisy speech spectra. For the multicondition training, the mean shift is not a problem.

HD removes harmonically related information, which is not perceptually important. This helps in the multicondition training experiment. The overall error rate reduction with HD for the multicondition training is 5.2%. Error reductions are 5.3%, 8.8% and 1.5% for Sets A, B, and C, respectively. The other techniques do not help in the multicondition training experiment. Half-wave rectification in the peak isolation algorithm appears to be harmful in the matched condition

because information on the valleys is removed. Peak-to-valley ratio locking also removes meaningful information on the height of the log spectral peak and hence is harmful to the matched training condition. VFR captures dynamic spectral changes and helps in the multicondition training.

Table 2 shows the results in multicondition training with HD, VFR and RASTA. The error rate reduction are: 11.52%, 22.34%, and 13.48% for Sets A, B, and C, respectively. The overall error reduction is 11.86%.

In summary, compared with the results in the previously submitted Eurospeech paper, there are two changes which improved ASR results. Voicing detection is added to speech classification, and VFR is applied on each speech segment instead of the whole utterance. By using the new speech detection technique, the overall error rate reduction for clean training is improved from 40.5% in the previous paper to 48.8% here using all the techniques without VFR, and to 53.01% with VFR. For the multicondition training experiment the overall error rate reduction is improved from 0.95% to 5.2% using only HD with the new speech detection technique, and to 11.86% by using HD, VFR and RASTA together.

## 6. References

- [1] H. G. Hirsch and D. Pearce, "The AURORA Experimental Framework for the Performance Evaluations of Speech Recognition Systems under Noisy Condition", *ISCA ITRW ASR2000 Automatic Speech Recognition: Challenges for the Next Millennium*, France, 2000.
- [2] B. Stroppe and A. Alwan, "A model of dynamic auditory perception and its application to robust word recognition," *IEEE Trans. on Speech and Audio Processing*, Vol. 5, No. 5, p. 451-464, 1997.
- [3] Q. Zhu and A. Alwan, "On the use of variable frame rate analysis in speech recognition", *Proc. IEEE ICASSP, Turkey*, Vol. III, p. 1783-1786, June 2000.

- [4] Q. Zhu and A. Alwan, "Amplitude Demodulation of Speech Spectra and its Application to Noise Robust Speech Recognition", *ICSLP 2000. Vol. 1*, pp. 341-344.
- [5] H. Hermansky and N. Morgan, "RASTA processing of speech". *IEEE Trans. Speech and Audio Proc.*, Vol.2, (no.4), Oct. 1994, p.578-89.
- [6] Medan, Y., Yair, E., and Chazan, D., "Super resolution pitch determination of speech signals". *IEEE Trans. on Signal Processing*, vol.39, (no.1), pp.40-8, Jan. 1991.