

Harmonic Structure Transform for Speaker Recognition

Kornel Laskowski & Qin Jin

Carnegie Mellon University, Pittsburgh PA, USA
KTH Speech Music & Hearing, Stockholm, Sweden

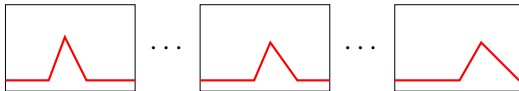
29 August, 2011

Spectral Transforms in General

Given $\mathbf{x} \equiv$ the energy spectrum of a speech frame,

$$\mathbf{y} = \mathcal{F}^{-1} \left(\log \left(\mathbf{M}^T \mathbf{x} \right) \right) - \langle \text{normalization term} \rangle$$

The matrix \mathbf{M} is a filterbank, whose columns look like:



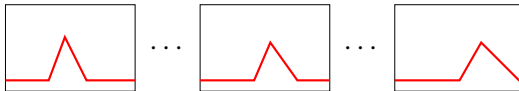
\mathbf{M} defines the **number** of filters, and their **central frequencies**, **widths**, and **general shapes**.

Spectral Transforms in General

Given $\mathbf{x} \equiv$ the energy spectrum of a speech frame,

$$\mathbf{y} = \mathcal{F}^{-1} \left(\log \left(\mathbf{M}^T \mathbf{x} \right) \right) - \langle \text{normalization term} \rangle$$

The matrix \mathbf{M} is a filterbank, whose columns look like:

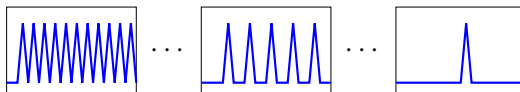


\mathbf{M} defines the **number** of filters, and their **central frequencies**, **widths**, and **general shapes**.

Importantly here, the filters of all such filterbanks integrate energy across frequencies **related by adjacency**.

The Harmonic Structure Transform (HST)

In contrast, the HST is implemented by a matrix \mathbf{H} whose columns look like:



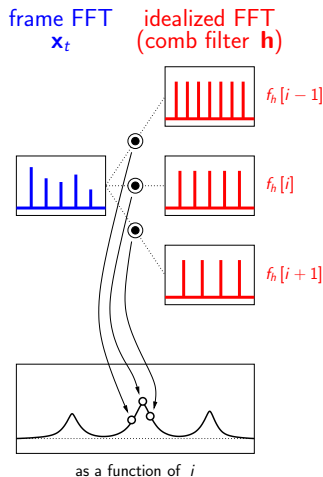
Each filter integrates energy across frequencies **related by harmonicity** (not adjacency).

- this is **novel** (Laskowski & Jin, 2010) for speaker recognition
- related to (Liénard, Barras & Signol, 2008) for pitch detection
- unknown: **number** of filters, and their **fundamental frequencies, tooth widths, and individual tooth shapes**

Outline of this Talk

- 1 Baseline Performance
 - What **is** known?
- 2 Experiments in HSCC Filterbank Design
 - linear spacing in fundamental frequency
 - piecewise linear spacing in fundamental frequency
 - logarithmic spacing in fundamental frequency
 - fundamental frequency range and density
- 3 Score-level Fusion with Standard MFCCs
- 4 Generalization
- 5 Conclusions

HST Processing



- analysis every 8 ms
- frames 32 ms wide
- comb filter teeth triangular (global width parameter)
- 400 filters, linearly spanning from 50 Hz to 450 Hz
- logarithm at each filter output, then normalization
- decorrelation using LDA
- yields *harmonic structure cepstral coefficients (HSCCs)*

HSCC Modeling for Classification

As simple as possible.

- one GMM per speaker
 - 1 assume one Gaussian element
 - 2 determine optimal number N_D of LDA dimensions
 - 3 hold N_D fixed
 - 4 determine optimal number of N_G Gaussians
- maximum likelihood closed-set classification (MAP under uniform prior)

Available Results (Laskowski & Jin, ODYSSEY 2010)

- Wall Street Journal data, mostly read speech
- 100-way closed-set classification, per gender
- ≈ 1500 10-second trials, per gender and dataset
- matched channel and matched multi-session conditions

System	Female, ♀		Male, ♂	
	DEV	TEST	DEV	TEST
<i>F0</i>	17.6	18.4	26.2	27.4
HST/LDA	99.7	99.9	99.7	99.7
MEL/DCT	98.7	99.3	99.3	98.6
MEL/LDA	98.7	99.3	99.3	98.9

Session Mismatch

- MIXER5 data, various speaking styles
- 66-way closed-set classification
- ≈ 3000 10-second trials, per dataset
- matched channel and matched session: accuracies of 100%
- matched channel but **mismatched session**:

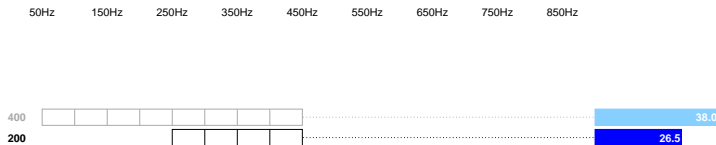
System	DEV	TEST
<i>F0</i>	14.1	16.2
HST/LDA	59.8	68.1
MEL/DCT	74.4	84.4
MEL/LDA	81.5	87.8

Linear Spacing of Fundamental Frequencies

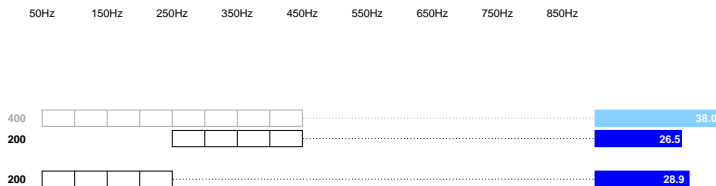
50Hz 150Hz 250Hz 350Hz 450Hz 550Hz 650Hz 750Hz 850Hz



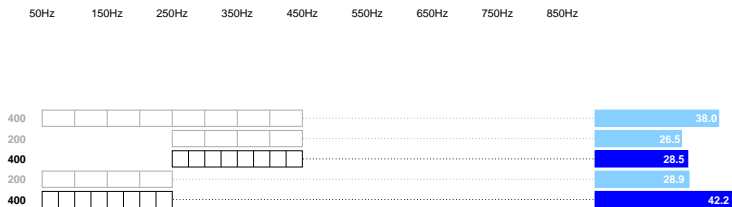
Linear Spacing of Fundamental Frequencies



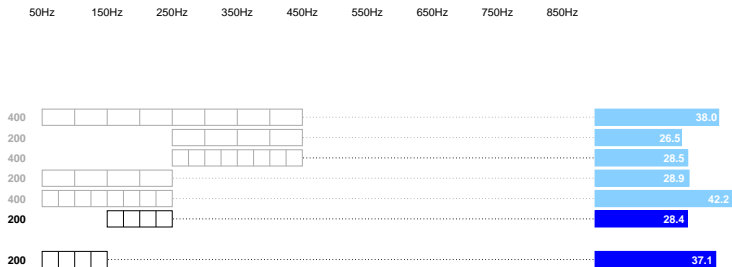
Linear Spacing of Fundamental Frequencies



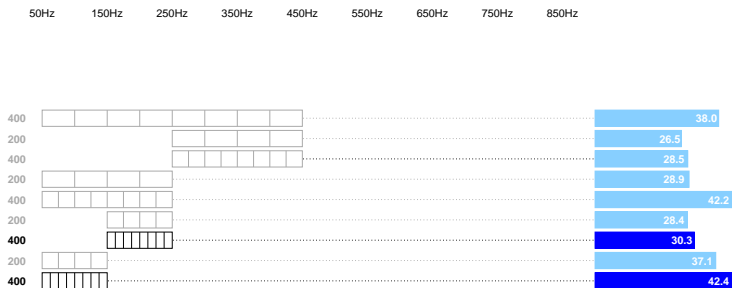
Linear Spacing of Fundamental Frequencies



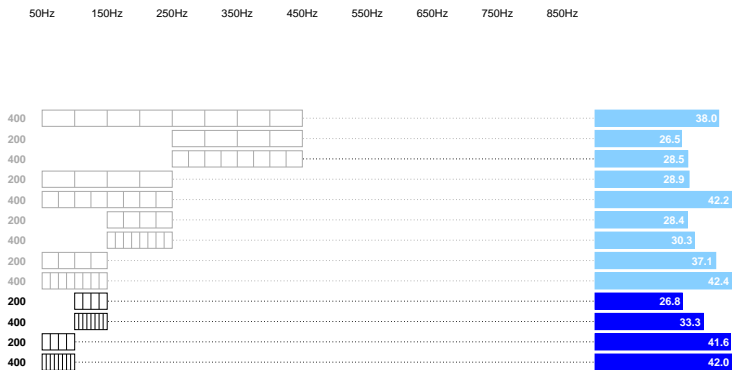
Linear Spacing of Fundamental Frequencies



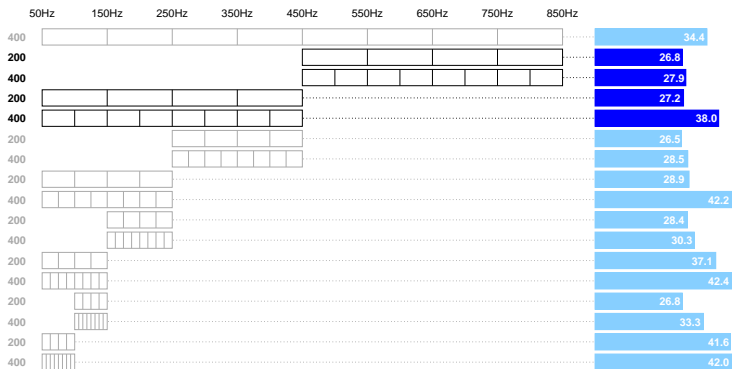
Linear Spacing of Fundamental Frequencies



Linear Spacing of Fundamental Frequencies



Linear Spacing of Fundamental Frequencies



Linear Spacing of Fundamental Frequencies

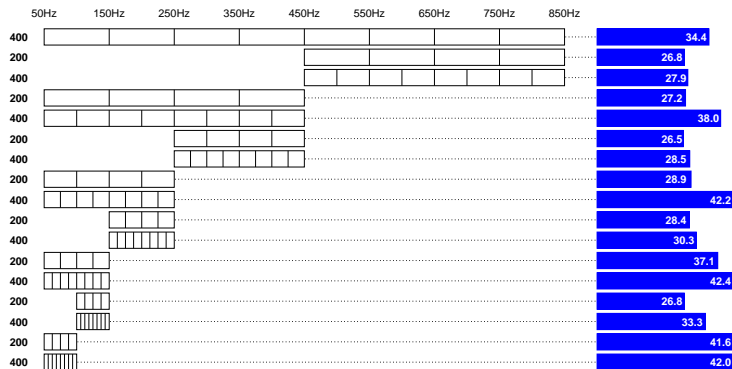
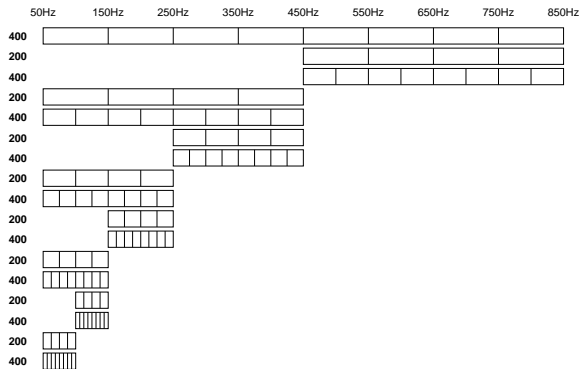
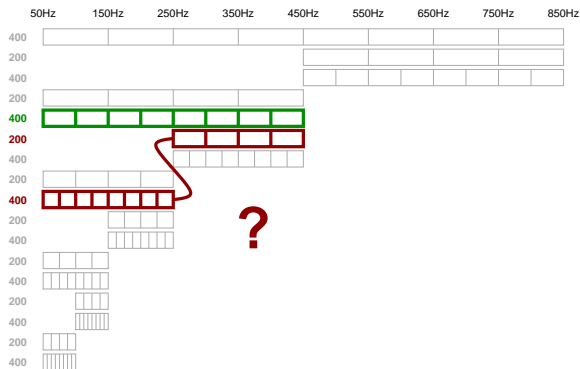


Figure 10 consists of two side-by-side frequency spectrum plots. The left plot shows the frequency spectrum of the speech signal, with frequency bins labeled from 50Hz to 850Hz. The right plot shows the frequency spectrum of the speech signal after processing, with the same frequency bins. The processed spectrum shows a significant reduction in the high-frequency components, particularly in the 400-850Hz range, which is highlighted by the red bars in the right plot. The red bars represent the magnitude of the high-frequency components, which are significantly lower than the blue bars representing the low-frequency components. The blue bars represent the magnitude of the low-frequency components, which are significantly higher than the red bars. The red bars are labeled with values ranging from 26.8 to 66.5, while the blue bars are labeled with values ranging from 26.8 to 66.5. The red bars are labeled with values ranging from 26.8 to 66.5, while the blue bars are labeled with values ranging from 26.8 to 66.5.

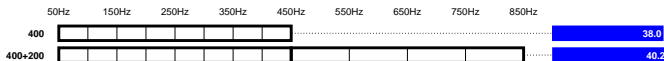
Linear Spacing of Fundamental Frequencies



Linear Spacing of Fundamental Frequencies



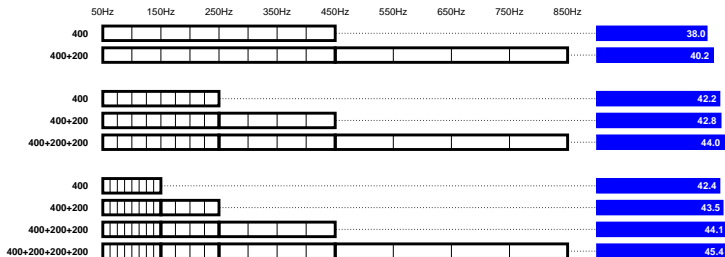
Piecewise Linear Spacing of Fundamental Frequencies



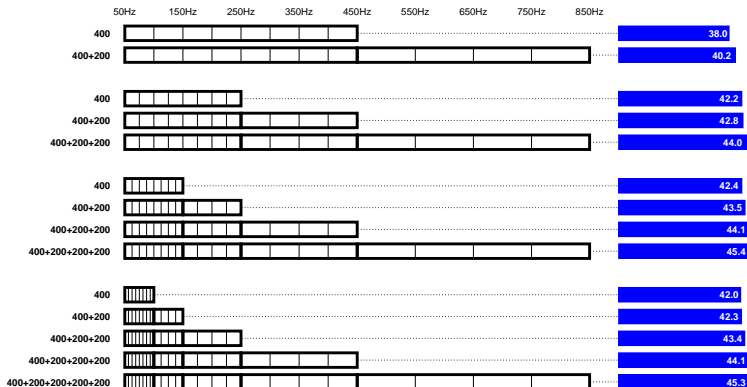
Piecewise Linear Spacing of Fundamental Frequencies



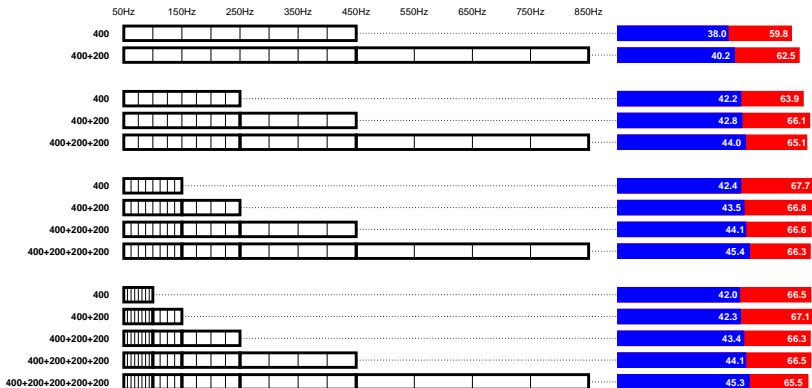
Piecewise Linear Spacing of Fundamental Frequencies



Piecewise Linear Spacing of Fundamental Frequencies



Piecewise Linear Spacing of Fundamental Frequencies

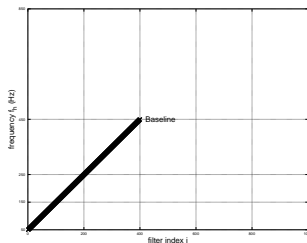


Logarithmic Spacing of Fundamental Frequencies

In the linear case:

- N_h fundamental frequencies f_h between f_{min} and f_{max} :

$$f_h[i] = f_h^{min} + \frac{i-1}{N_h-1} (f_h^{max} - f_h^{min}) \quad 1 \leq i \leq N_h$$



Baseline

38.0

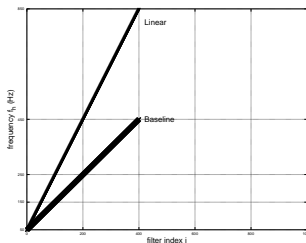
59.8

Logarithmic Spacing of Fundamental Frequencies

In the linear case:

- N_h fundamental frequencies f_h between f_{min} and f_{max} :

$$f_h[i] = f_h^{min} + \frac{i-1}{N_h-1} (f_h^{max} - f_h^{min}) \quad 1 \leq i \leq N_h$$



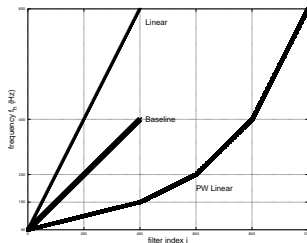
Baseline	38.0	59.8
Linear	34.4	56.7

Logarithmic Spacing of Fundamental Frequencies

In the linear case:

- N_h fundamental frequencies f_h between f_{min} and f_{max} :

$$f_h[i] = f_h^{min} + \frac{i-1}{N_h-1} (f_h^{max} - f_h^{min}) \quad 1 \leq i \leq N_h$$



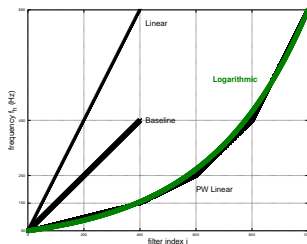
Baseline	38.0	59.8
Linear	34.4	56.7
PW Linear	45.4	66.3

Logarithmic Spacing of Fundamental Frequencies

In the linear case:

- N_h fundamental frequencies f_h between f_{min} and f_{max} :

$$f_h[i] = f_h^{min} + \frac{i-1}{N_h-1} (f_h^{max} - f_h^{min}) \quad 1 \leq i \leq N_h$$



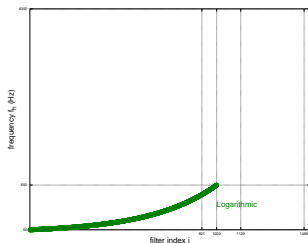
Baseline	38.0	59.8
Linear	34.4	56.7
PW Linear	45.4	66.3
Logarithmic	45.3	66.0

$$f_h[i] = f_h^{min} \left(\frac{f_h^{max}}{f_h^{min}} \right)^{\frac{i-1}{N_h-1}}$$

Range and Density of Fundamental Frequencies

Three manipulations:

- f_{min} : raise to 62.5 Hz
- f_{max} : raise to maximize accuracy (to 4000 Hz)
- N_h : lower to maximize accuracy (to 1129)

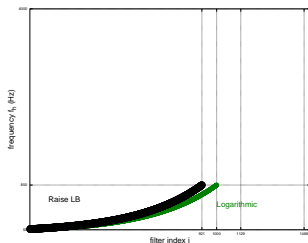


Baseline	38.0	59.8
Linear	34.4	56.7
PW Linear	45.4	66.3
Logarithmic	45.3	66.0

Range and Density of Fundamental Frequencies

Three manipulations:

- f_{min} : raise to 62.5 Hz
- f_{max} : raise to maximize accuracy (to 4000 Hz)
- N_h : lower to maximize accuracy (to 1129)

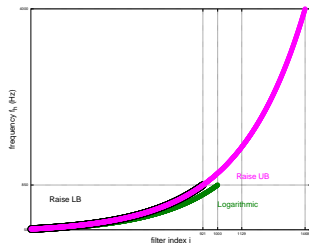


Baseline	38.0	59.8
Linear	34.4	56.7
PW Linear	45.4	66.3
Logarithmic	45.3	66.0
Raise LB	43.4	64.7

Range and Density of Fundamental Frequencies

Three manipulations:

- f_{min} : raise to 62.5 Hz
- f_{max} : raise to maximize accuracy (to 4000 Hz)
- N_h : lower to maximize accuracy (to 1129)

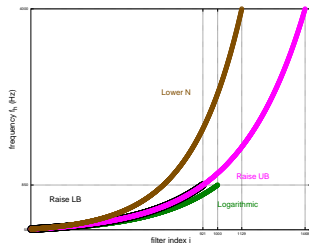


Baseline	38.0	59.8
Linear	34.4	56.7
PW Linear	45.4	66.3
Logarithmic	45.3	66.0
Raise LB	43.4	64.7
Raise UB	50.6	70.2

Range and Density of Fundamental Frequencies

Three manipulations:

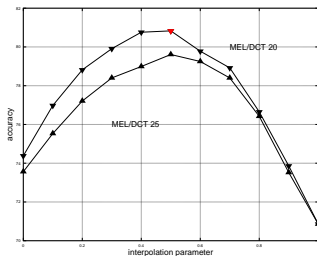
- f_{min} : raise to 62.5 Hz
- f_{max} : raise to maximize accuracy (to 4000 Hz)
- N_h : lower to maximize accuracy (to 1129)



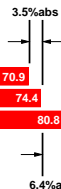
Baseline	38.0	59.8
Linear	34.4	56.7
PW Linear	45.4	66.3
Logarithmic	45.3	66.0
Raise LB	43.4	64.7
Raise UB	50.6	70.2
Lower N	51.2	70.9

Interpolation with Standard MFCCs

- linear interpolation with DCT-decorrelated log-Mel energies
 - 20 coefficients
 - 25 coefficients (always worse)

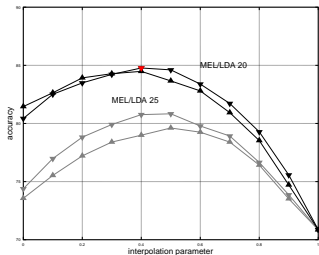


HSCC
MEL/DCT
+ HSCC



Interpolation with LDA-Rotated MFCCs

- linear interpolation with LDA-decorrelated log-Mel energies
 - 20 coefficients (better in combination)
 - 25 coefficients (better alone)

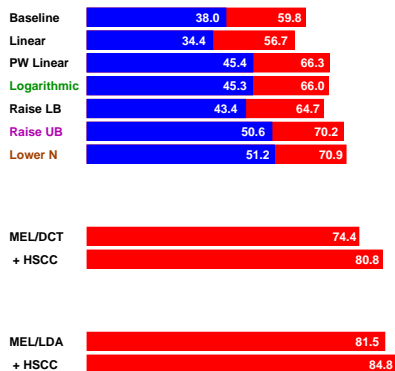


HSCC
MEL/LDA
+ HSCC



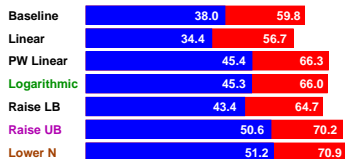
Summary of Accuracy on DEVSET

DEVSET

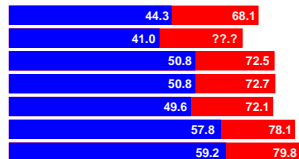


Accuracy on EVALSET

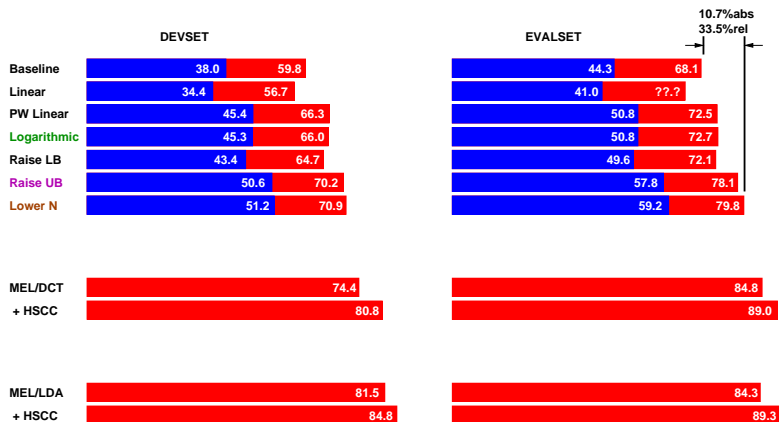
DEVSET



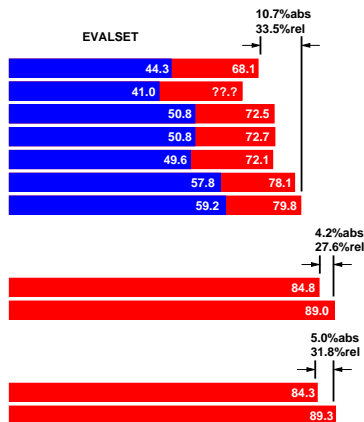
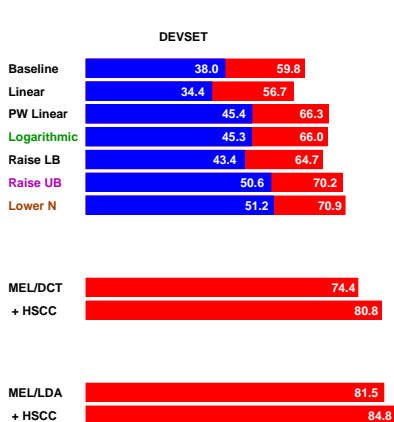
EVALSET



Accuracy on EVALSET



Accuracy on EVALSET



Conclusions

- ① evaluated the baseline transform in **session mismatch**
 - viable: twice as many errors an equivalent MFCC system
- ② errors can be reduced **by a third (33.5%rel)** by optimizing:
 - the number of filters in the filterbank
 - the fundamental frequency corresponding to each filter
- ③ **logarithmic spacing** of fundamental frequencies is better than linear spacing
 - more filters for low fundamental frequencies
 - fewer filters for high fundamental frequencies
- ④ in an equivalent MFCC system, errors can be reduced **by almost a third (27.6-31.8%rel)** via score-level fusion with the improved HST system

Future Directions

- ❶ change **framing policy** from 8 ms/32 ms to something longer
 - intonation and voice quality are **supra**-segmental
 - larger temporal support → greater spectral resolution
- ❷ optimize **the tooth shape** of comb filters
- ❸ find a **data-independent** decorrelation transform
 - leading to a compact (< 25 coefficients) representation
- ❹ explore adaptation from a **universal background model**
- ❺ generalize to **binary speaker verification** (and NIST SREs)

Potential Impact

- ① a new **general** representation of the spectrum
- ② deliberately **orthogonal** to spectral envelope features (MFCCs, LPCCs, etc.)
 - but computed in an identical manner
- ③ likely beneficial not only for speaker recognition, but also:
 - online **speaker diarization**
 - classification of “emotional speech”
 - clinical voice quality assessment

Some Interesting Insights ...

- ① currently in HST, spectral energy < 300 Hz is zeroed out
 - this improves closed-set speaker classification
 - but **the fundamental (zeroth harmonic) is ignored**
 - the fundamental is thought to play a role in **emotional expression**
- ② that optimal filter spacing is logarithmic is curious
 - independent of the logarithmic tonotopicity of the basilar membrane
 - greater acuity in discriminating among harmonic sounds (not just pure tones)

THANK YOU