MODELING INSTANTANEOUS INTONATION FOR SPEAKER IDENTIFICATION USING THE FUNDAMENTAL FREQUENCY VARIATION SPECTRUM

Kornel Laskowski and Qin Jin

interACT, Carnegie Mellon University, Pittsburgh PA, USA Cognitive Systems Lab, Universität Karlsruhe, Karlsruhe, Germany

Goal

Improve speaker identification by modeling **intonation bias** (the distribution of speaker-preferred pitch change in octaves per 8 ms).

HOW?

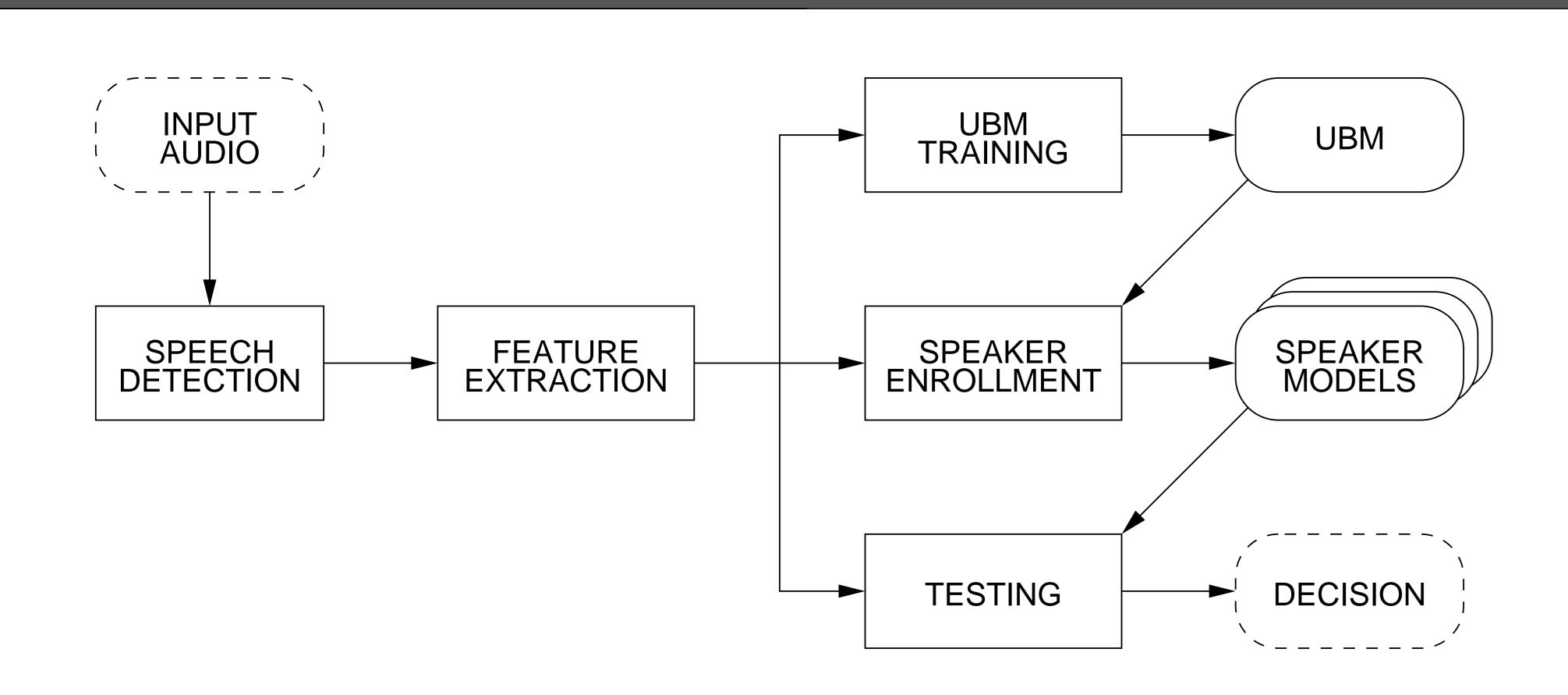
Compute the fundamental frequency variation (FFV) spectrum and model alongside "traditional" spectral (MFCC) features, ie.

$$\log P \text{ (MFCC13} \oplus \text{FFV7} | \mathcal{M} \text{)} = \log P \text{ (MFCC13, FFV7} | \mathcal{M}_{stacked} \text{)}$$

or

 $\log P \left(\text{MFCC13} \otimes \text{FFV7} \middle| \mathcal{M} \right) = \lambda_{\text{MFCC13}} \log P \left(\text{MFCC13} \middle| \mathcal{M}_{\text{MFCC13}} \right)$ $+ \lambda_{\text{FFV7}} \log P \left(\text{FFV7} \middle| \mathcal{M}_{\text{FFV7}} \right)$

Baseline GMM-MFCC SID System



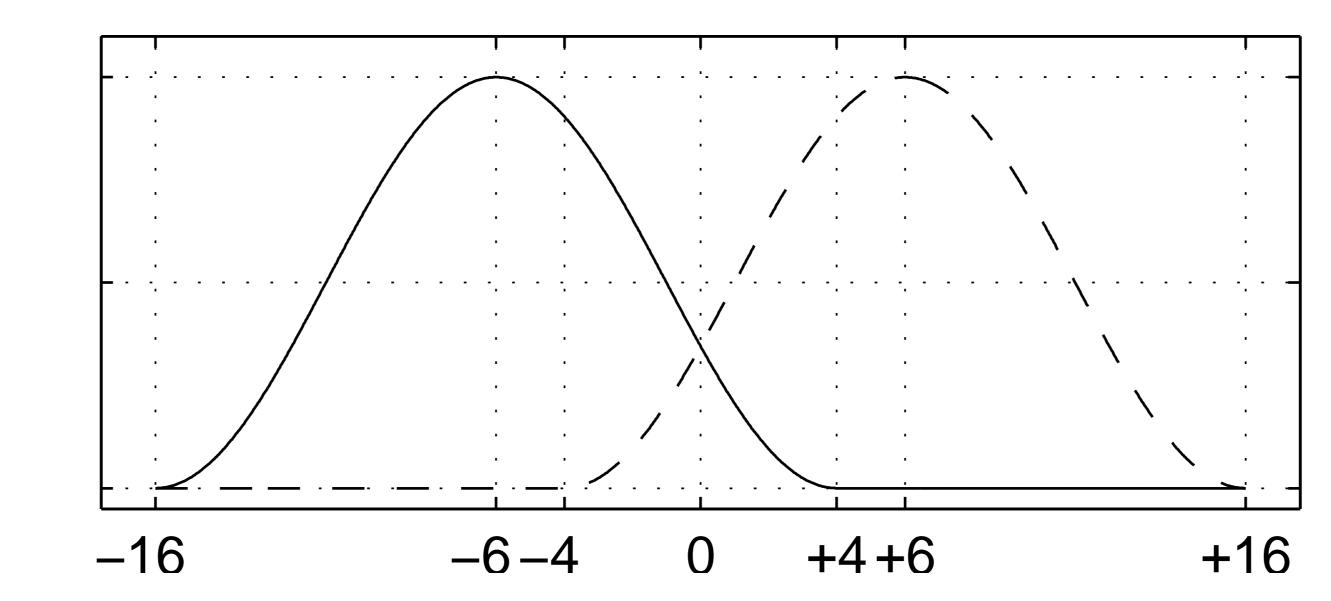
- 32 ms frames every 8 ms, low-energy frames excluded
- MFCC13 $\equiv \{MFCC_0, MFCC_1, \cdots, MFCC_{12}\}$
- cepstral mean subtraction
- Gaussian mixture models (4096)
- ML estimation of universal background model (UBM)
- MAP estimation of speaker GMM means

Data

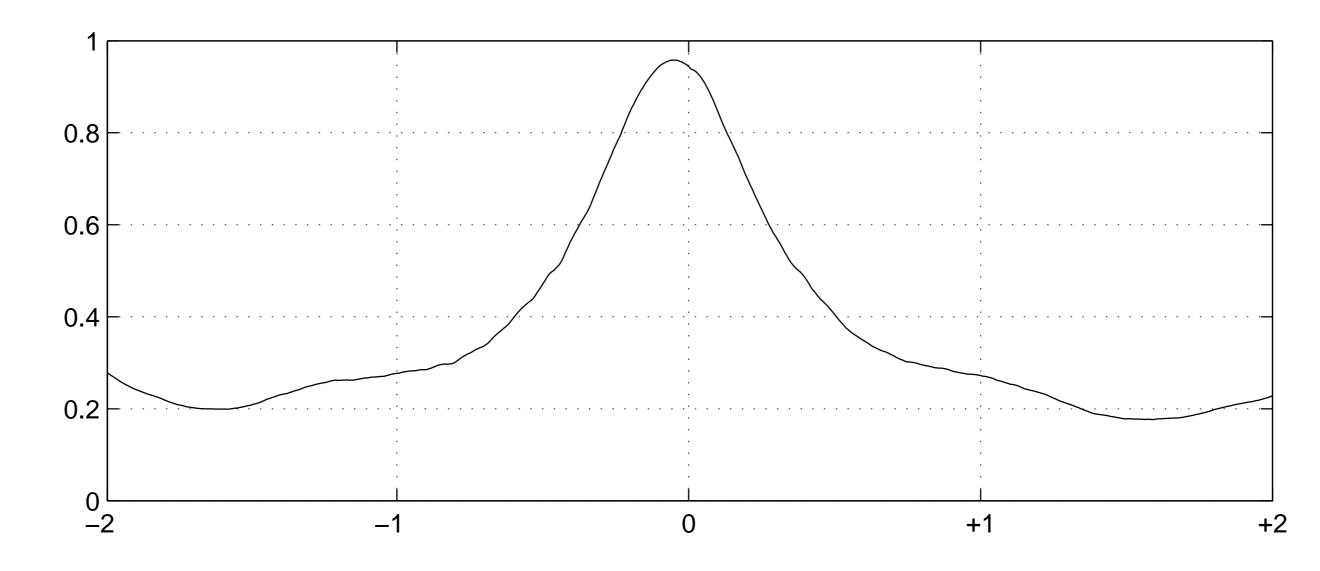
- Wall Street Journal, LDC CSR-I (WSJ0) & LDC CSR-II (WSJ1)
- read sentences and some spontaneously produced utterances
- 16 kHz wideband audio (Sennheiser HMD414)
- 102 female speakers & 95 male speakers
 TRAINSET: 5 minutes of speech per speaker
 TESTSET: 3 trials × 1 minute of speech per speaker
- UBMSET: remaining speakers' speech (70 hours)

What Is the Fundamental Frequency Variation Spectrum?

• dot product of two spectra, F_L and F_R ...



- ... after dilating one of \mathbf{F}_L and \mathbf{F}_R by factor $\mathbf{2}^{\rho}$
- ullet FFV spectrum is value of dot product over range of ho



• after 7-filter filterbank, $FFV7 \equiv \{FFV_{-3}, \cdots, FFV_{+3}\}$

Findings

- 1. Raw FFv7 features offer better than chance performance.
- 2. Sphering FFv7 features dramatically improves performance.
- 3. Model-space combination of MFCC13 with sphered FFV7 yields 54% and 40% relative error reductions for female and male speakers, respectively.
- 4. Model-space combinations of MFCC13 with FFV7 relative to that with MFCC7 yields 24–35% and 26–30% relative error reductions for female and male speakers, respectively.
- 5. Feature-space combinations of MFCC13 with FFV7 offer no benefit.

Experiments (identification accuracies, %)

System	Female	Male
MFCC13	82.0	92.3
MFCC7	43.8	68.8
FFV7	27.8	45.3
PCA(MFCC13)	84.3	91.2
PCA(MFCC7)	44.4	65.2
PCA(FFV7)	62.7	64.2
MFCC13 \otimes MFCC7	86.3	92.6
MFCC13 \otimes FFV7	80.7	92.3
PCA(MFCC13) ⊗ MFCC7	87.3	92.3
PCA(Mfcc13) ⊗ Ffv7	84.6	92.6
MFCC13 \otimes PCA(MFCC7)	89.2	93.7
MFCC13 \otimes PCA(FFV7)	91.8	95.4
PCA(MFCC13) ⊗ PCA(MFCC7)	86.9	93.0
PCA(MFCC13) ⊗ PCA(FFV7)	91.5	95.1