

Learning Prosodic Sequences Using the Fundamental Frequency Variation Spectrum

Problem

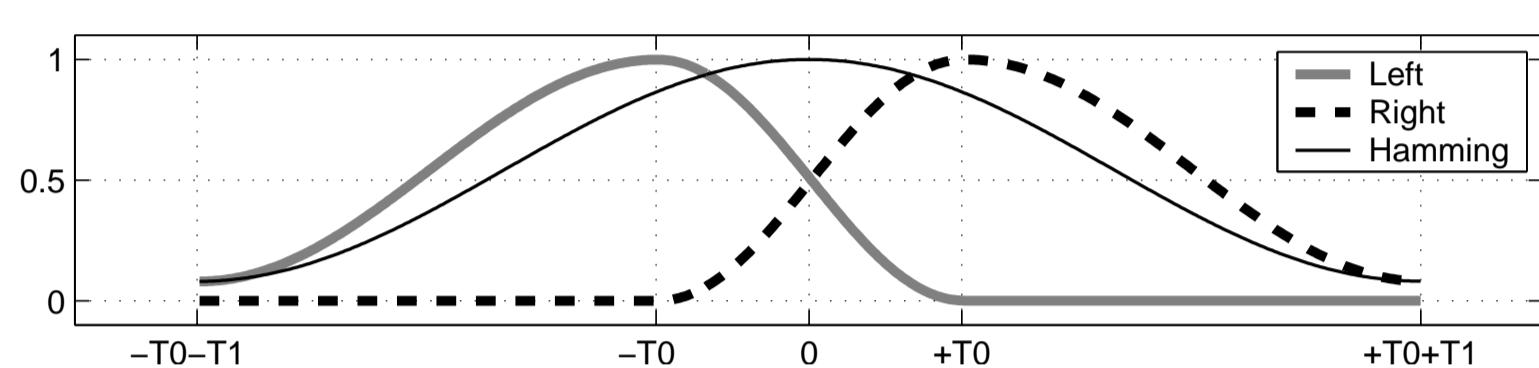
Current state-of-the-art conversational spoken dialogue systems are not sufficiently responsive. They produce speech at detected end-of-utterance (EOU) locations; EOU detection consists of **waiting for 0.5–2.0 seconds**.

Goal

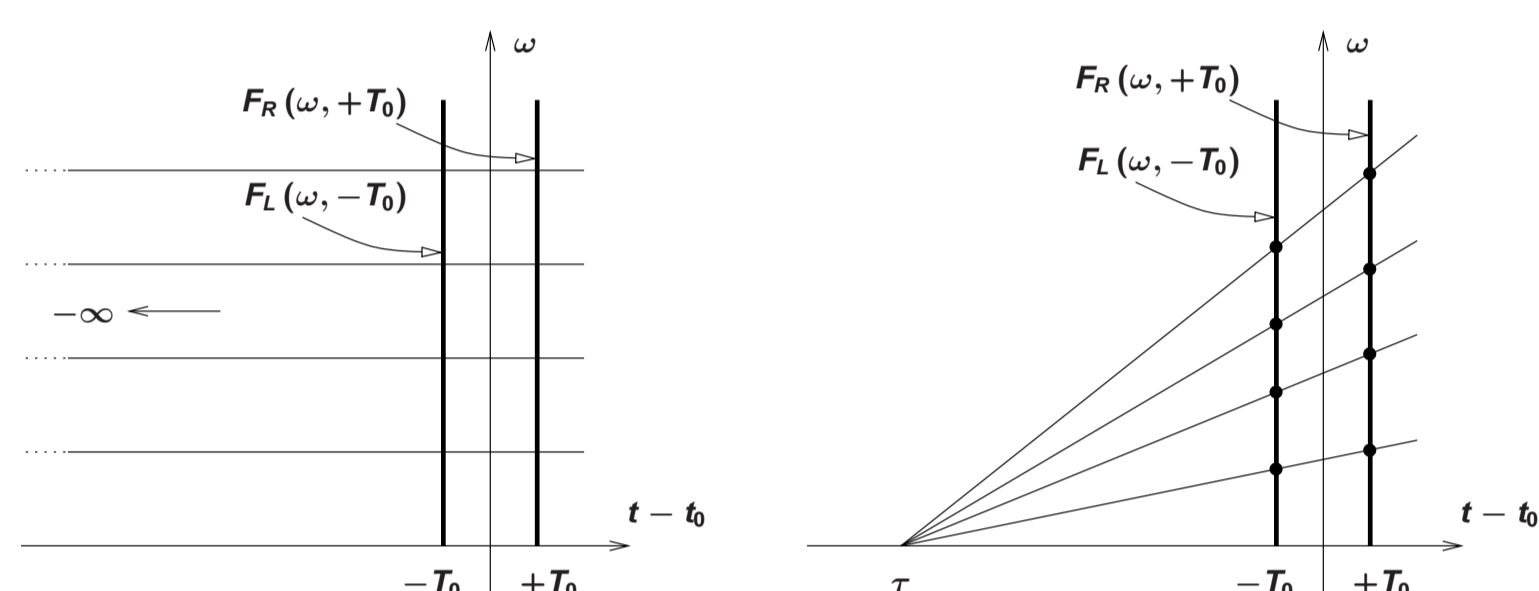
Faster online prediction (≤ 0.3 s) of end-of-utterance (EOU) locations.

The Fundamental Frequency Variation Spectrum

- ▶ use **entire** spectrum to quantify variation in F0
- ▶ sample spectrum at two locations in each frame: $-T_0$ and $+T_0$ relative to midpoint of frame



- ▶ define the **vanishing-point product** at a vanishing point τ



$$g^T(\tau) = \begin{cases} \int_{-f_s/2}^{+f_s/2} F_L\left(\frac{-\tau-T_0}{-\tau+T_0}f\right) F_R^*(f) df, & \tau < -T_0 \\ \int_{-f_s/2}^{+f_s/2} F_L(f) F_R^*\left(\frac{+\tau-T_0}{+\tau+T_0}f\right) df, & \tau > +T_0 \end{cases}$$

- ▶ define the conformal mapping

$$\rho = \begin{cases} -\log_2\left(\frac{-\tau-T_0}{-\tau+T_0}\right), & \tau < -T_0 \\ +\log_2\left(\frac{+\tau-T_0}{+\tau+T_0}\right), & \tau > +T_0 \end{cases}$$

and transform $g^T(\tau)$ to yield

$$g^\rho(\rho) = \begin{cases} \int_{-f_s/2}^{+f_s/2} F_L(f) F_R^*(2^\rho f) df, & \rho < 0 \\ \int_{-f_s/2}^{+f_s/2} F_L(2^{-\rho} f) F_R^*(f) df, & \rho \geq 0 \end{cases}$$

- ▶ define the linear interpolation

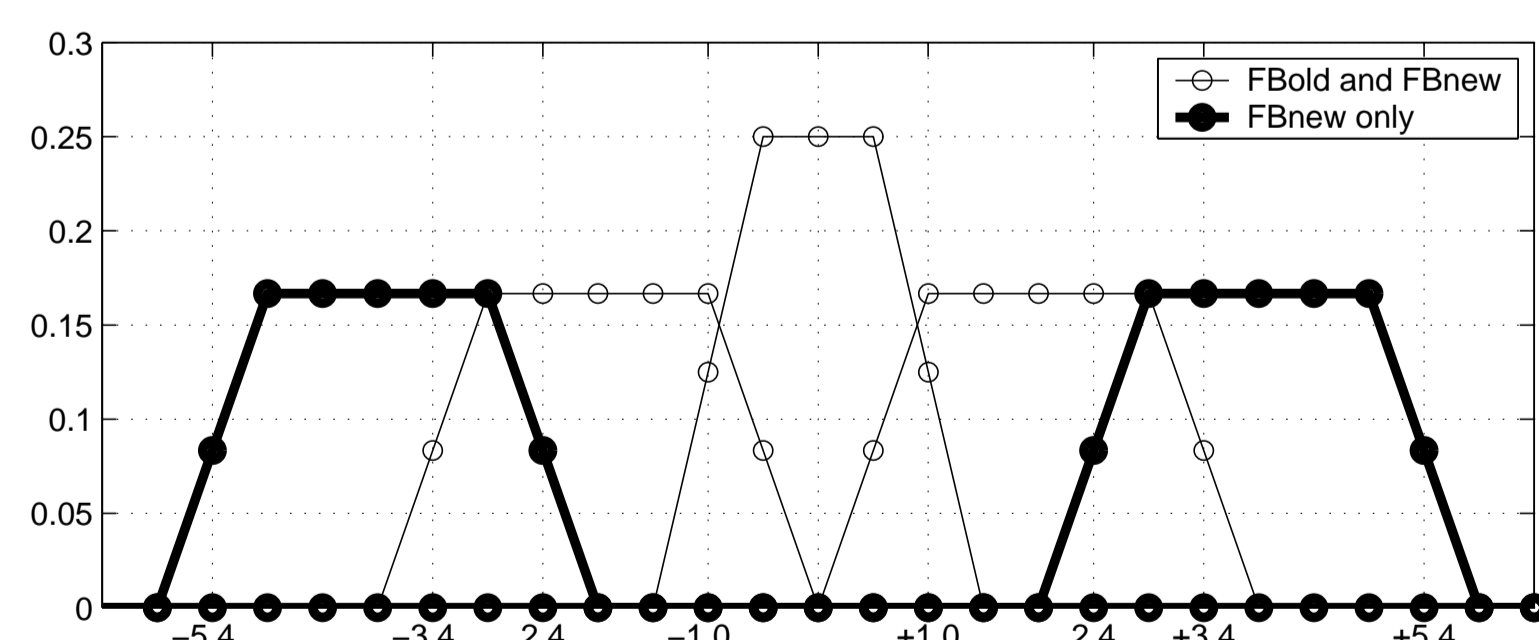
$$|\tilde{F}(2^{\pm\rho}k)| = \beta |F[2^{\pm\rho}k]| + (1-\beta) |F[2^{\pm\rho}k]|$$

where $\beta = |2^{\pm\rho}k| - 2^{\pm\rho}k|$

- ▶ sample $g^\rho(\rho)$ at discrete locations to yield

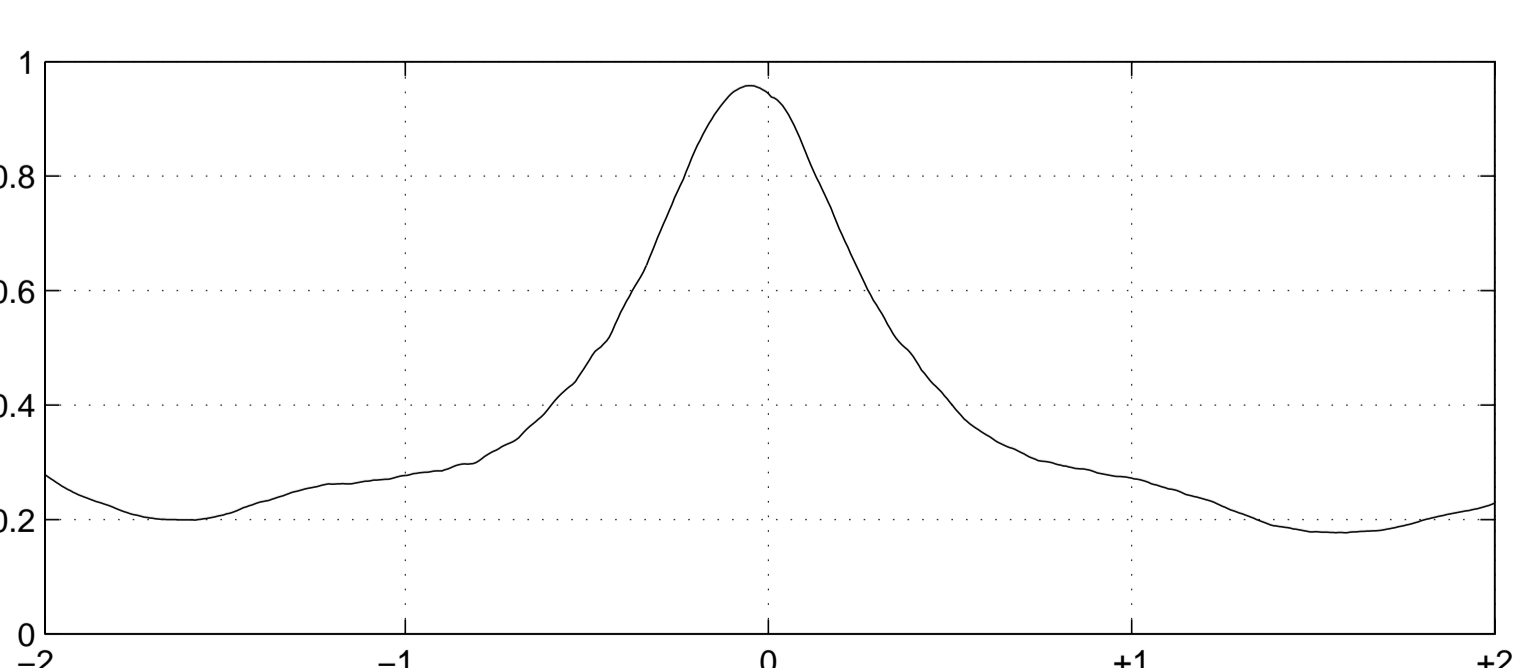
$$g^\rho[r] = \begin{cases} \sum_{k=-N/2+1}^{N/2} |\tilde{F}_L(2^{-4r/N}k)| |F_R^*[k]|, & r \geq 0 \\ \sum_{k=-N/2+1}^{N/2} |F_L[k]| |\tilde{F}_R^*(2^{4r/N}k)|, & r < 0 \end{cases}$$

- ▶ normalize for energy-independence, and apply filterbank

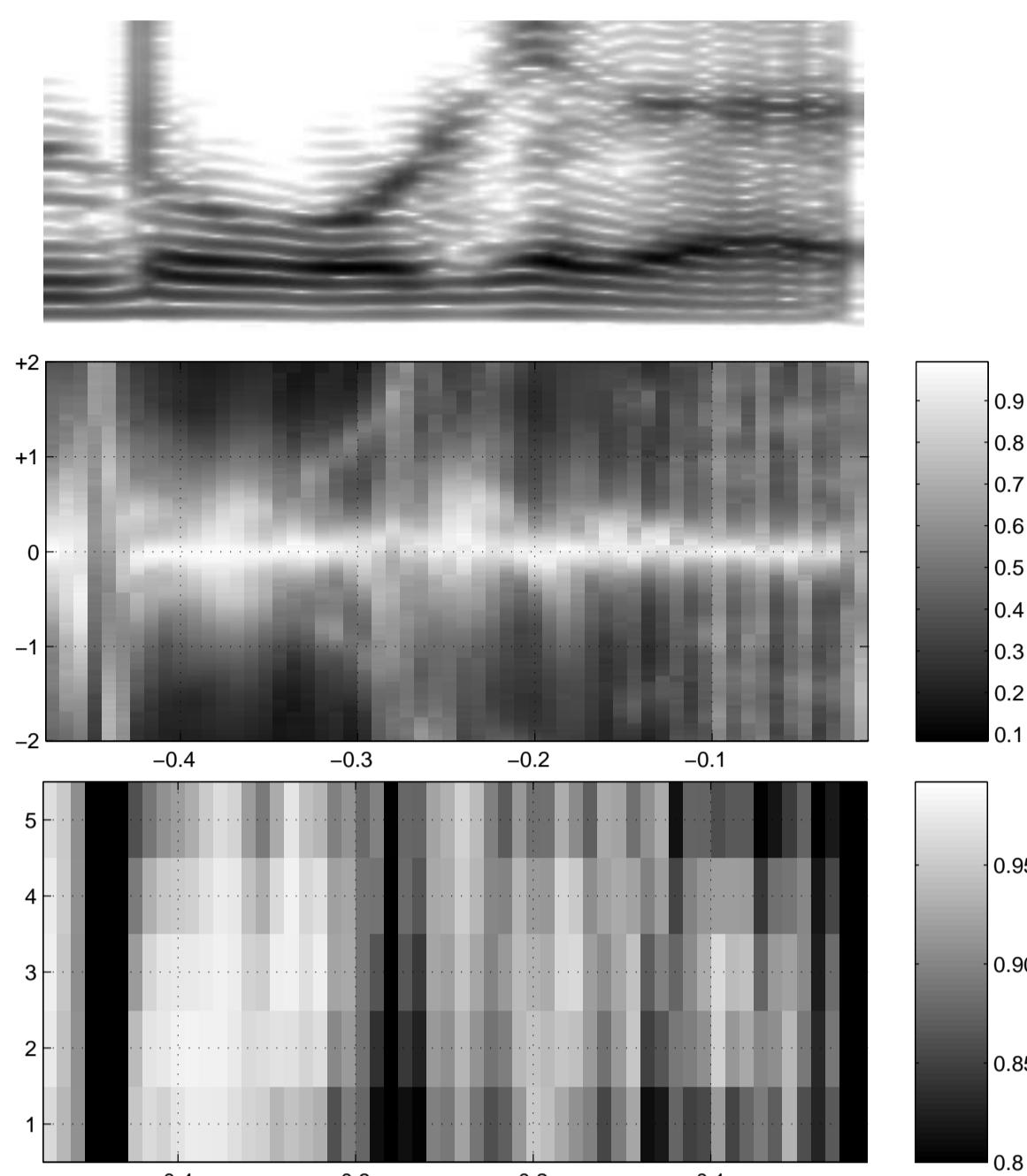


- ▶ apply Karhunen-Loève whitening transform

Sample "Spectrogram" Representation



Sample "Spectrogram" Representation



Modeling

Hidden Markov model for each of SC and \neg SC:

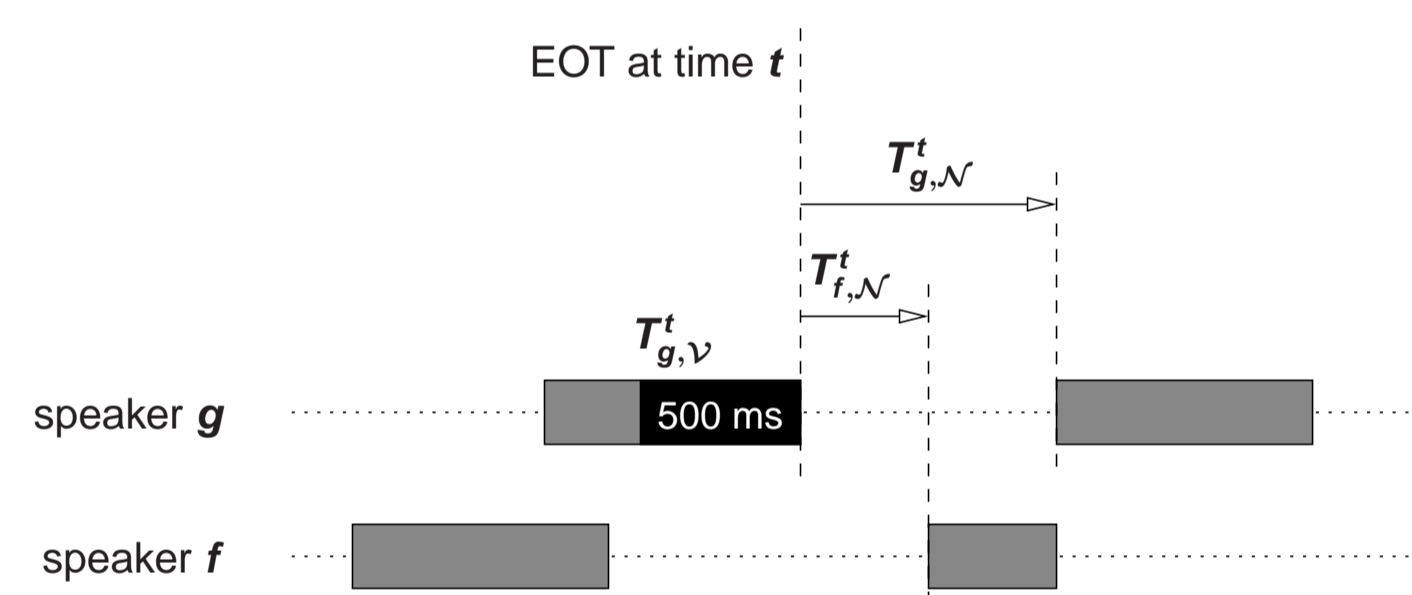
- ▶ 4 states
- ▶ 1 Gaussian per state
- ▶ trained on DEVSET using the Forward-Backward Algorithm

Online Speaker Change Prediction

Binary classification of each end-of-talkspurt (EOT) location in a human-human dialogue as either

- ▶ a *speaker change*, SC; or as
- ▶ *not a speaker change*, \neg SC

using a **prosodic description of 500 ms of audio** preceding the EOT.



Reference labels are given by the **automatic assignment**:

$$L_t = \begin{cases} \text{SC} & \text{if } T_{f,N}^t - T_{g,N}^t < 0 \\ \neg\text{SC}, & \text{otherwise} \end{cases}$$

In previous work, we have demonstrated that the occurrence of *observed* speaker changes at time t is strongly correlated with human judgement that they are *appropriate*.

Data

Swedish Map Task Corpus:

- ▶ two speakers: a *giver*, g , and a *follower*, f
- ▶ task: g explains directions to f

Data Set	Duration (mn:ss)	Dialogue role g		
		speakers	# EOTs	# SCs
DEVSET	77:40	F4,F5,M2,M3	480	222
EVALSET	60:39	F1,F2,F3,M1	317	149

- ▶ highly interactive dialogues
- ▶ DEVSET and EVALSET are disjoint in speakers

Automatic Classification

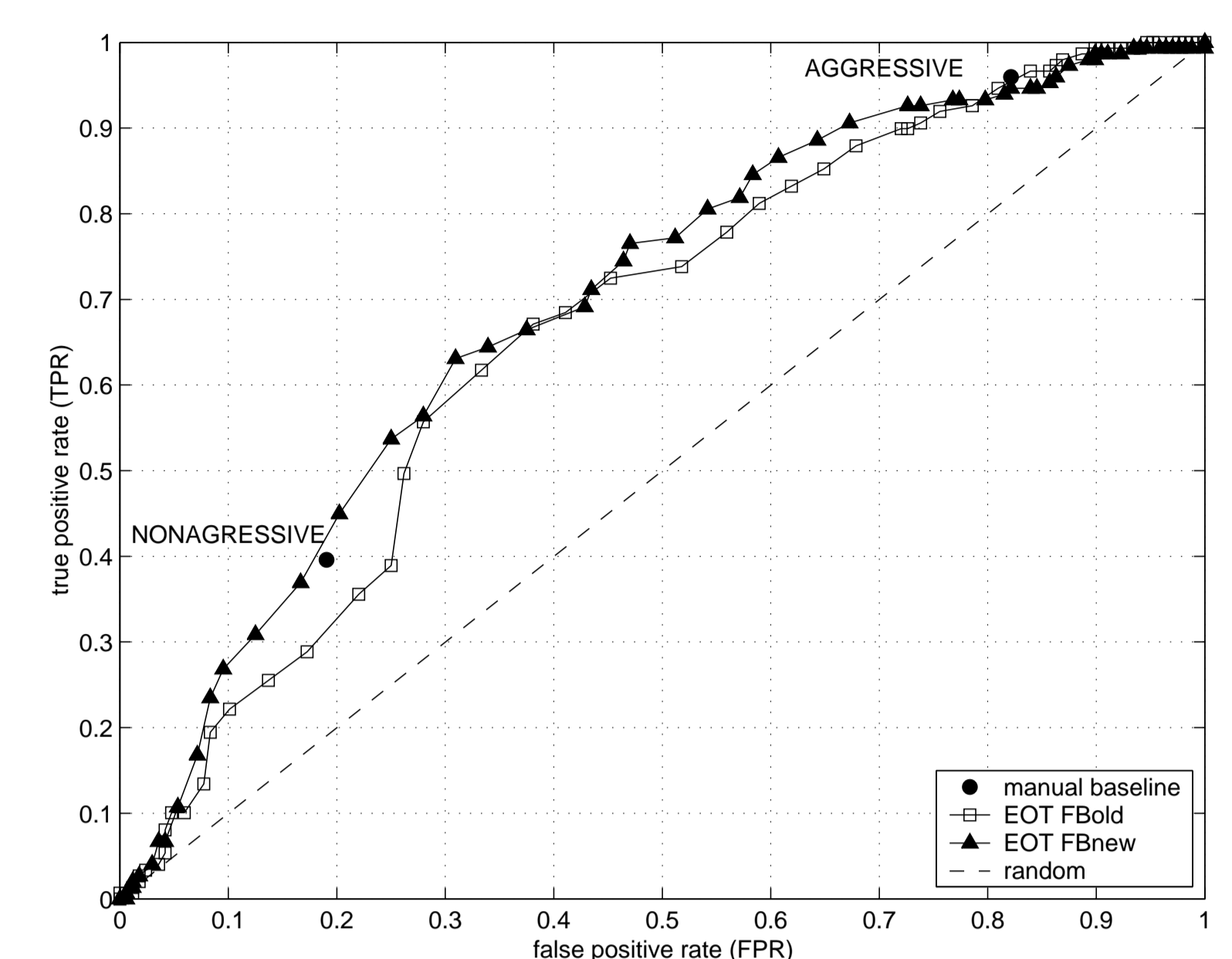
- ▶ compare two different log-likelihood-ratio classifiers
- ▶ train 10 HMMs for each of SC and \neg SC, by using different random seeds prior to Forward-Backward training
- ▶ **classifier 1**: log-likelihood-ratio over **mean** of 10 model likelihoods

$$L_t = \arg \max_k P(x | \mathcal{M}_k) \quad (1)$$

- ▶ **classifier 2**: log-likelihood-ratio over **product** of 10 model likelihoods

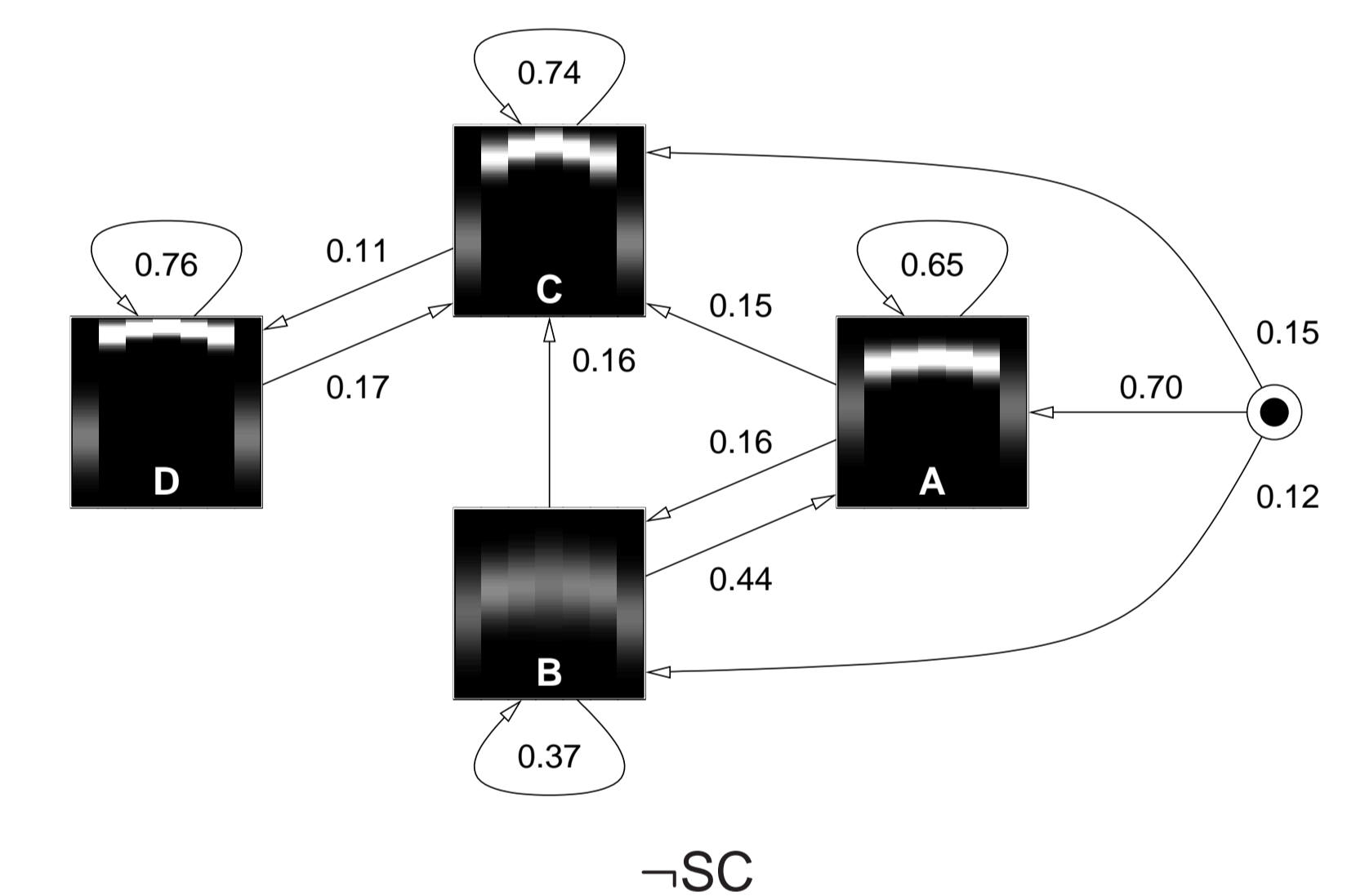
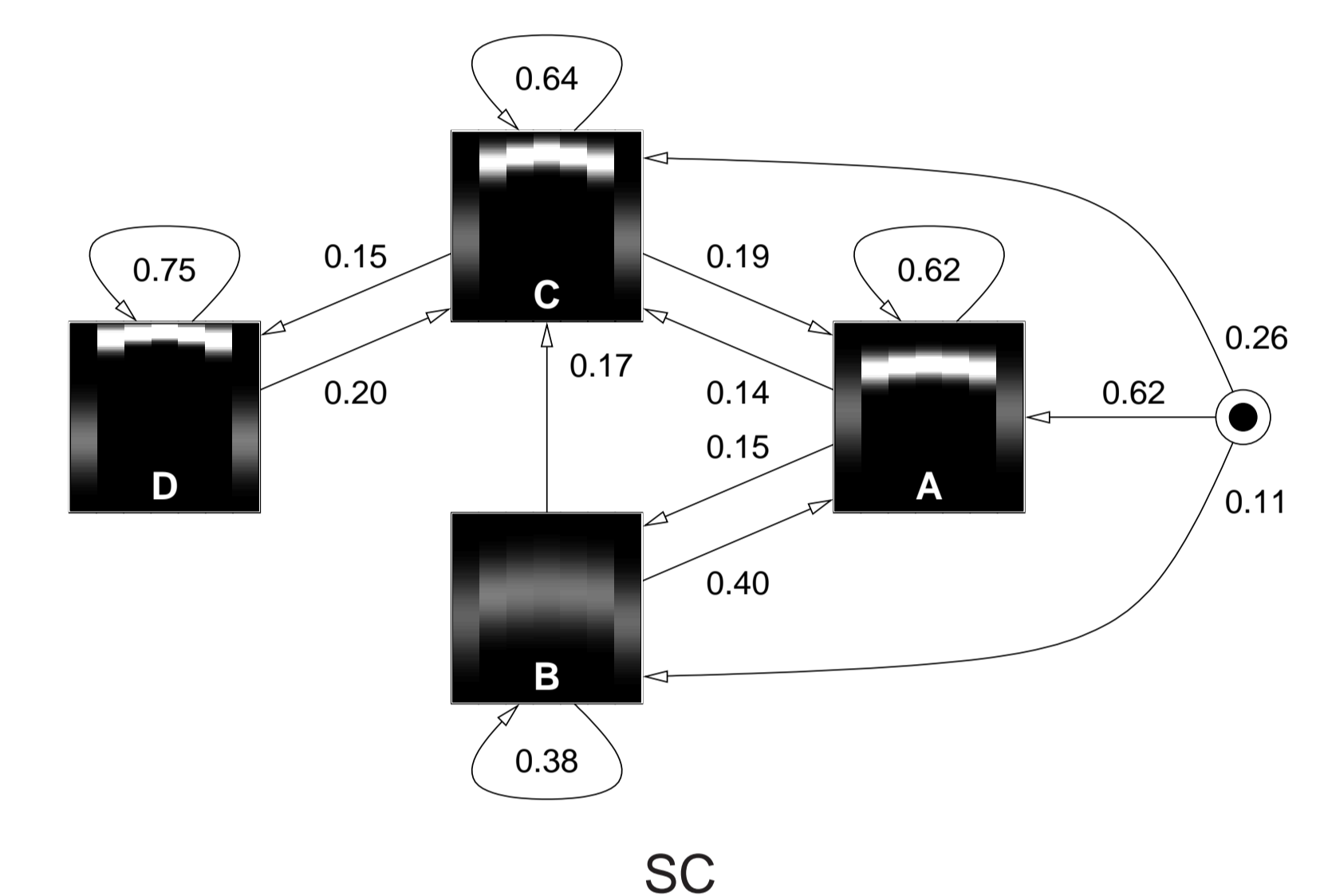
$$L_t = \arg \max_k \prod_{i=1}^{10} P(x | \mathcal{M}_{k,i}) \quad (2)$$

Experimental Results



Learned Sequences

- ▶ randomly chosen HMMs, showing transition and (unnormalized) emission probabilities:



- ▶ normalized emission probabilities for state C in both models:



Conclusions

- ▶ first exploration of what models of fundamental frequency variation sequences actually learn
- ▶ learned models corroborate existing research of human behavior
- ▶ representation appears suitable for direct, principled, continuous sequence modeling as in SAD and ASR
- ▶ improved filterbank design yields SC/ \neg SC classification accuracy improvements on unseen data of 12-17% relative