# Speech Separation for Recognition and Enhancement

Dan Ellis

Laboratory for Recognition and Organization of Speech and Audio
Dept. Electrical Eng., Columbia University, NY
*and*
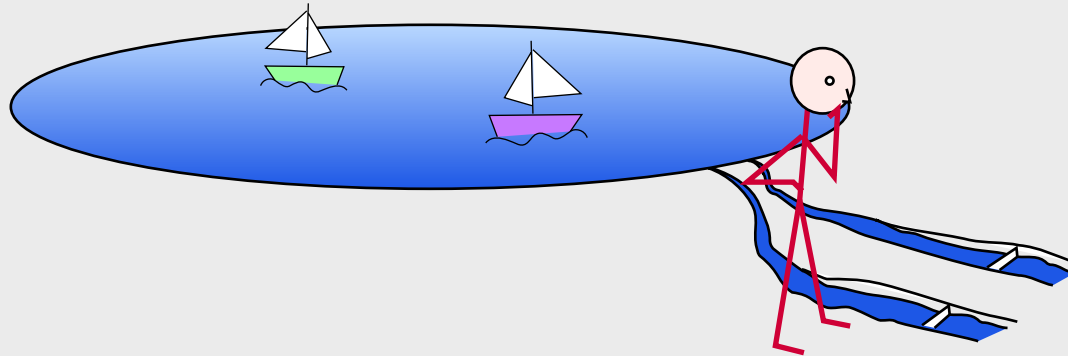International Computer Science Institute, Berkeley CA

dpwe@ee.columbia.edu          http://labrosa.ee.columbia.edu/

Lab
ROSA
Laboratory for the Recognition and
Organization of Speech and Audio

ICSI

# 1. Speech in the Wild

- **The world is cluttered**
  **sound is transparent**
  - mixtures are inevitable

- **Useful information is structured by 'sources'**
  - specific definition of a 'source':
    intentional independence

# Speech in the Wild: Examples

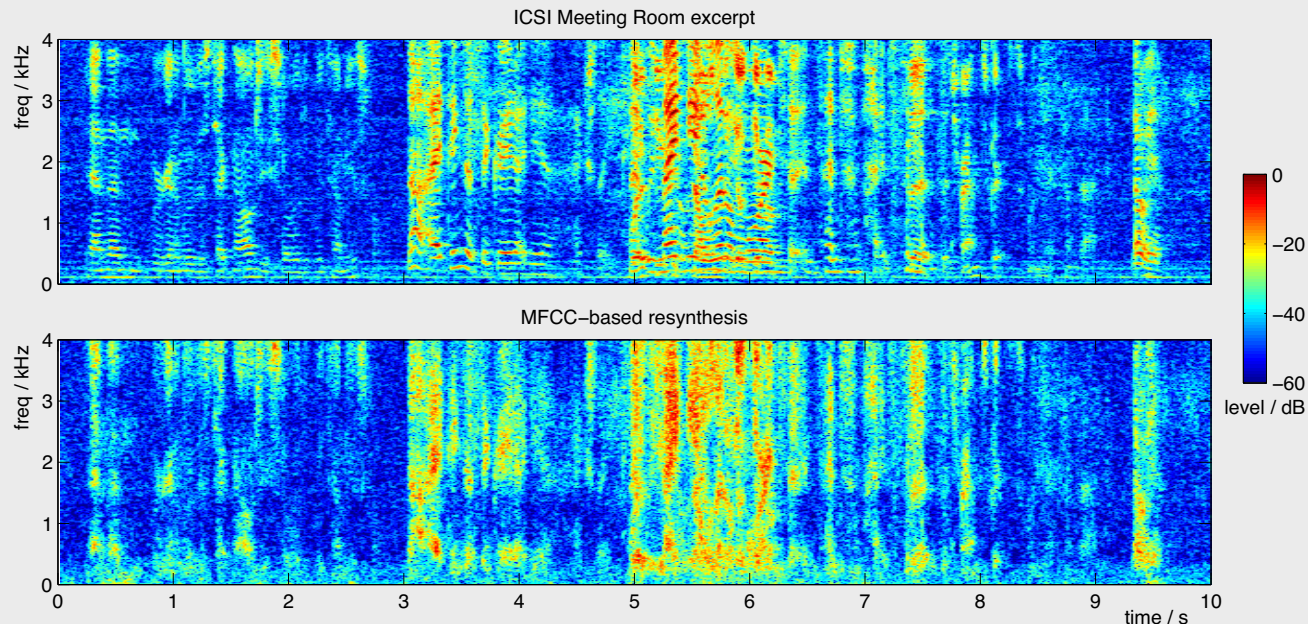- Multi-party discussions



- Ambient recordings



- Applications:
  - communications
  - robots
  - lifelogging/archives

# Recognizing Speech in the Wild

- Current ASR relies on low-D representations
  - e.g. 13 dimensional MFCC features every 10ms

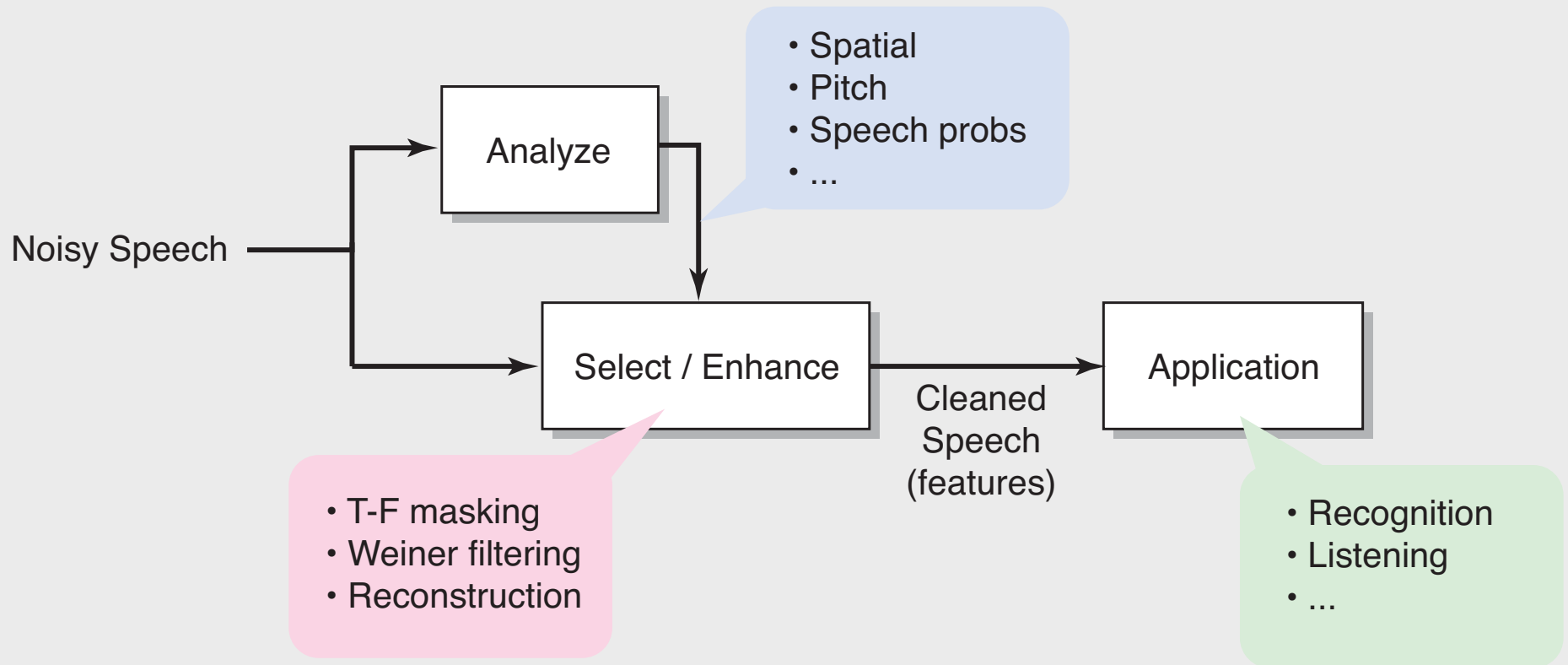ICSI Meeting Room excerpt

MFCC−based resynthesis

- very successful for clean speech!
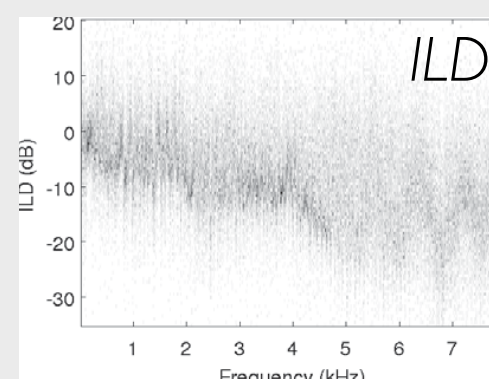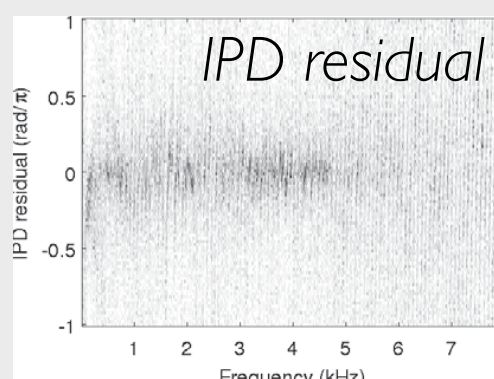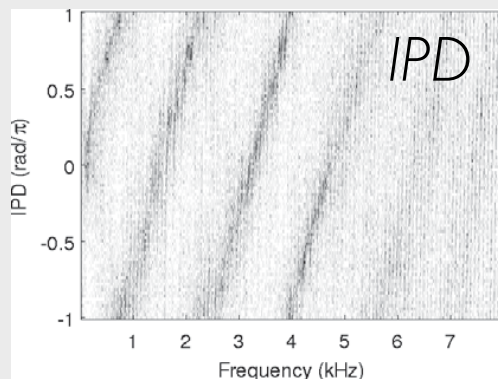- inadequate for mixtures

- We need separation!

# Separation by Spatial Info

- Given multiple microphones,
  sound carries spatial information about source
- E.g. model interaural spectrum of each source
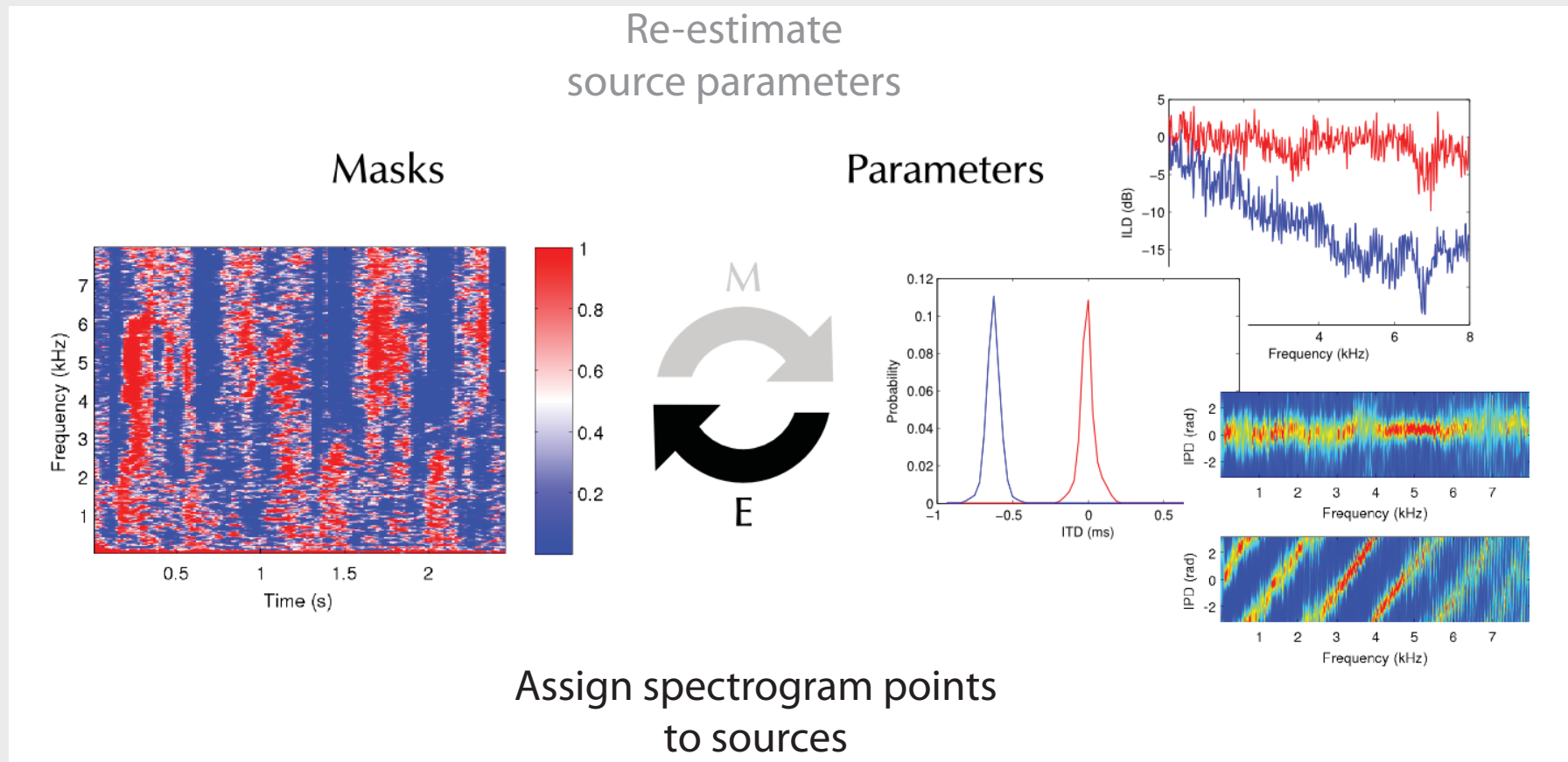  as stationary level and time differences:

$$\frac{L(\omega, t)}{R(\omega, t)} = a(\omega)e^{j\omega\tau}N(\omega, t)$$

- e.g. at 75°, in reverb:

# Model-Based EM Source Separation and Localization (MESSL)
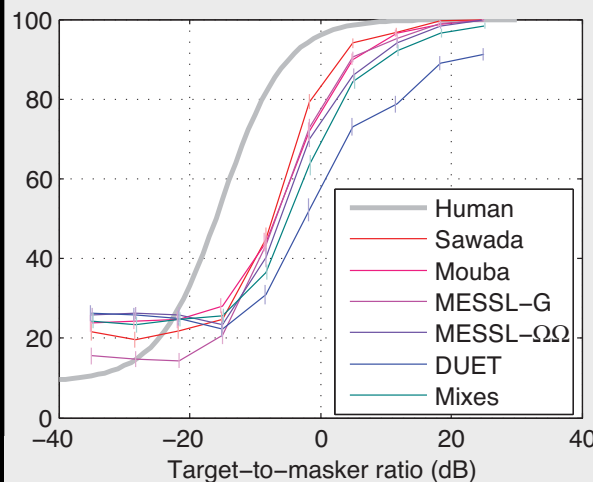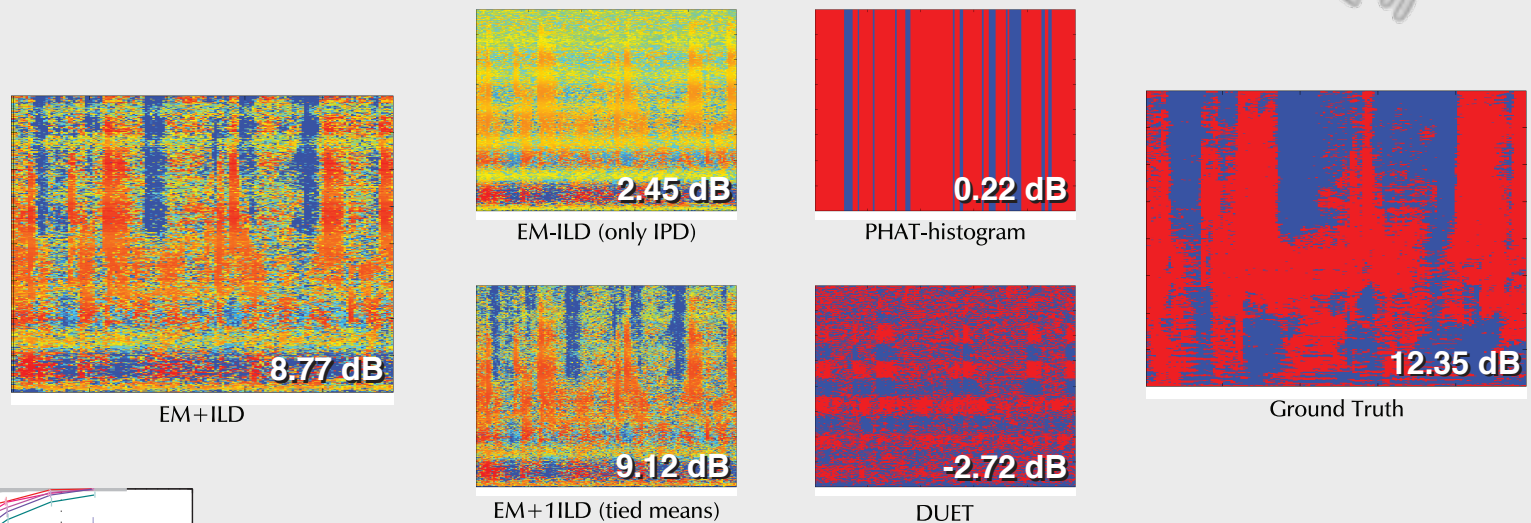
*Mandel et al. '10*



Re-estimate source parameters

Assign spectrogram points to sources

○ can model more sources than sensors

# MESSL Results

- ## Modeling uncertainty improves results
  - tradeoff between constraints & noisiness


EM+ILD — 8.77 dB


EM-ILD (only IPD) — 2.45 dB


PHAT-histogram — 0.22 dB


EM+1ILD (tied means) — 9.12 dB
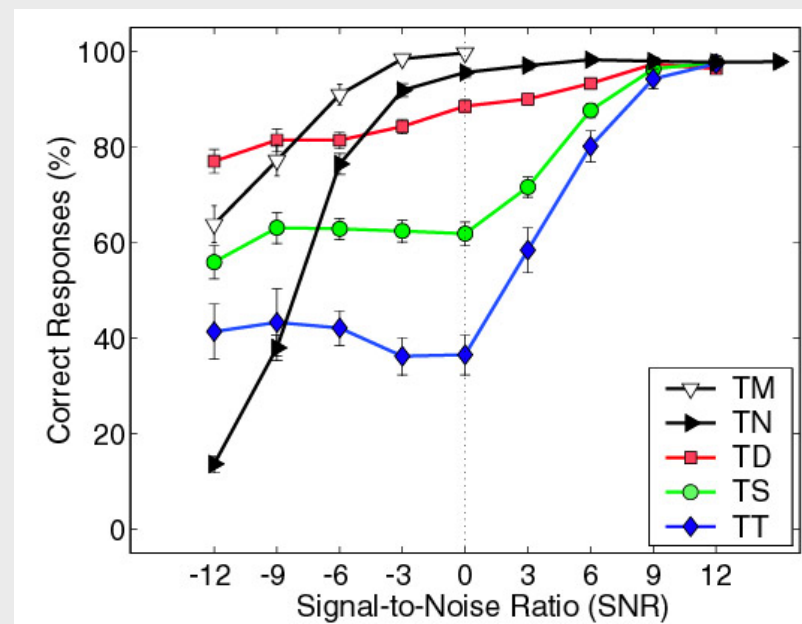

DUET — -2.72 dB


Ground Truth — 12.35 dB



- ## Helps with recognition
  - digits accuracy

# 3. Separation by Pitch

- Voiced syllables have near-periodic "pitch"

  - perceptually salient
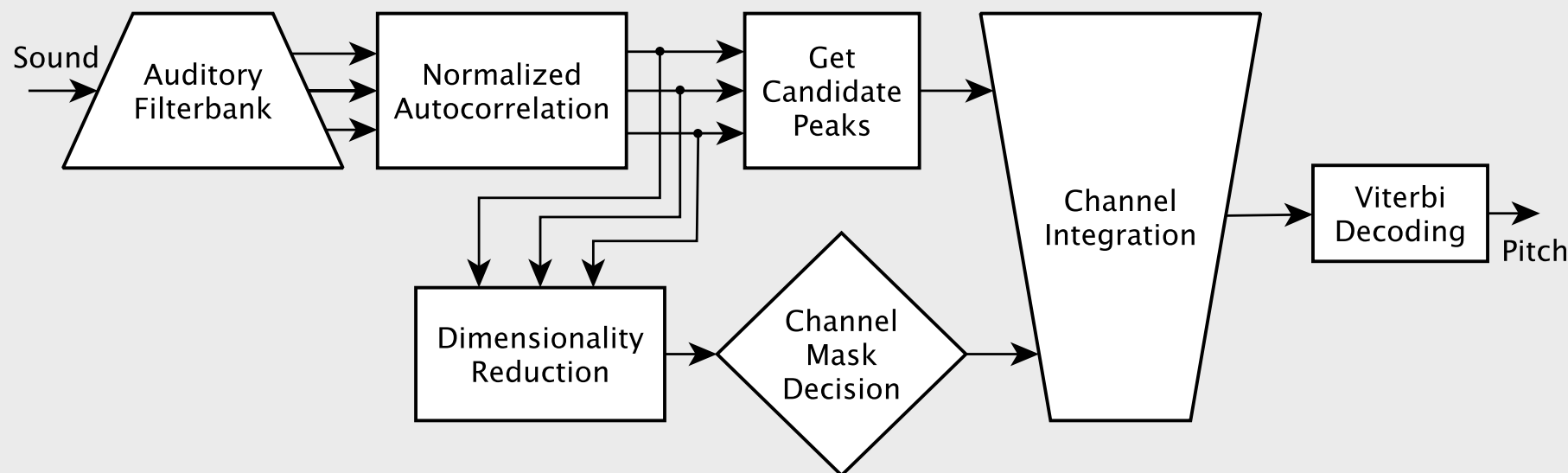
  - lost in MFCCs



*Brungart et al.'01*

- Can we track pitch & use it for separation?
  - ... and other speech tasks?
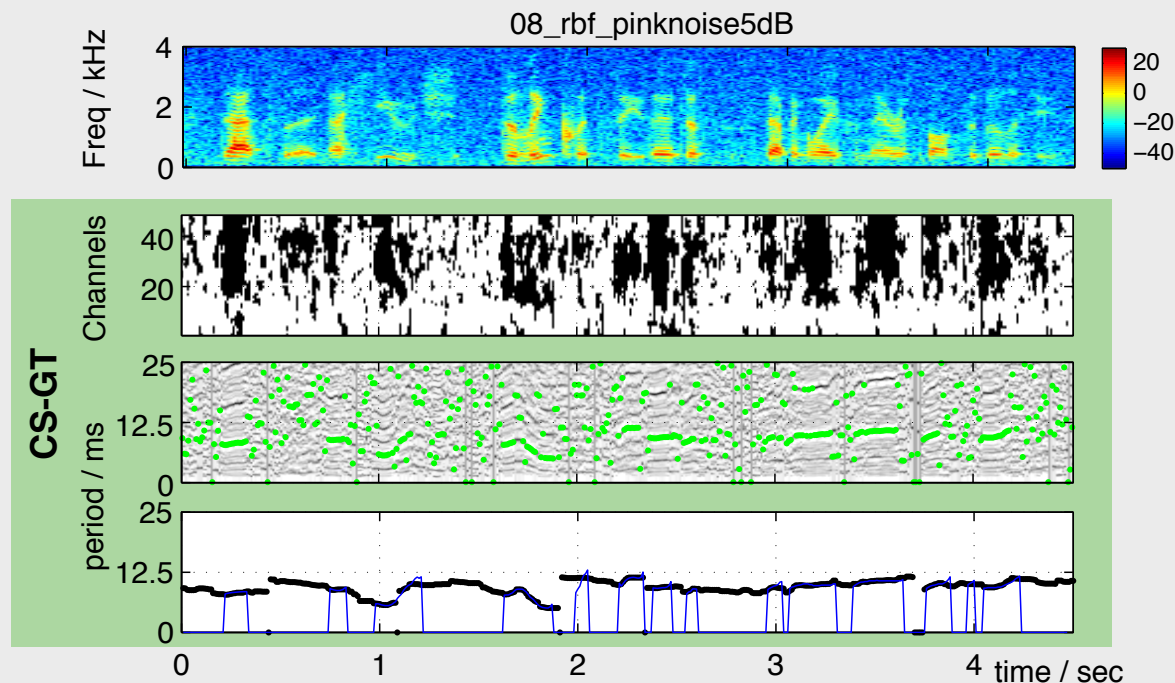
# Noise-Robust Pitch Tracking

*BS Lee & Ellis '12*

- **Important for voice detection & separation**
- **Based on** channel selection *Wu, Wang & Brown '03*
  - ○ pitch from summary autocorrelation over "good" bands



  - ○ trained classifier decides which channels to include

# Noise-Robust Pitch Tracking

- **Channel-based classifiers**
  learn domain channel/noise characteristics
  - then separate, or derive features for recognition
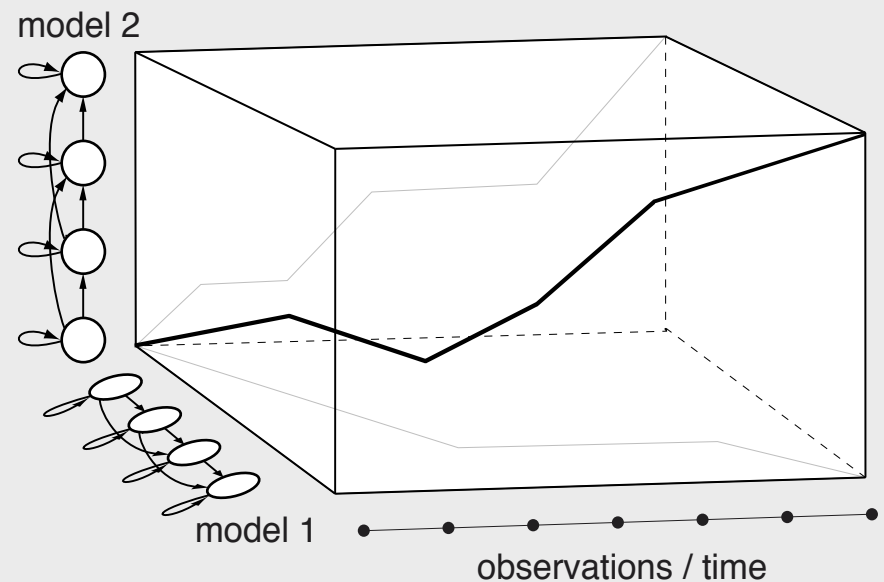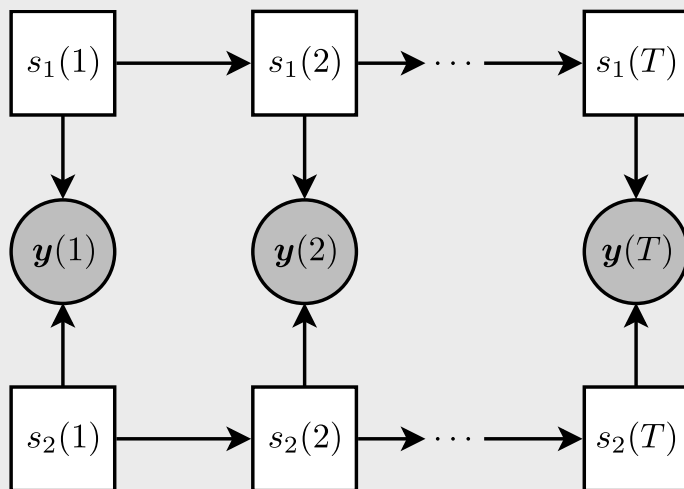


08_rbf_pinknoise5dB

- Only works for pitched sounds
  - need a broader description of the speech source...

# 4. Separation by Models

*Varga & Moore, '90*
*Hershey et al., '10*

- If ASR is finding best-fit parameters
  argmax P($W$|$X$) ...

- Recognize mixtures with Factorial HMM
  - model + state sequence for each voice/source
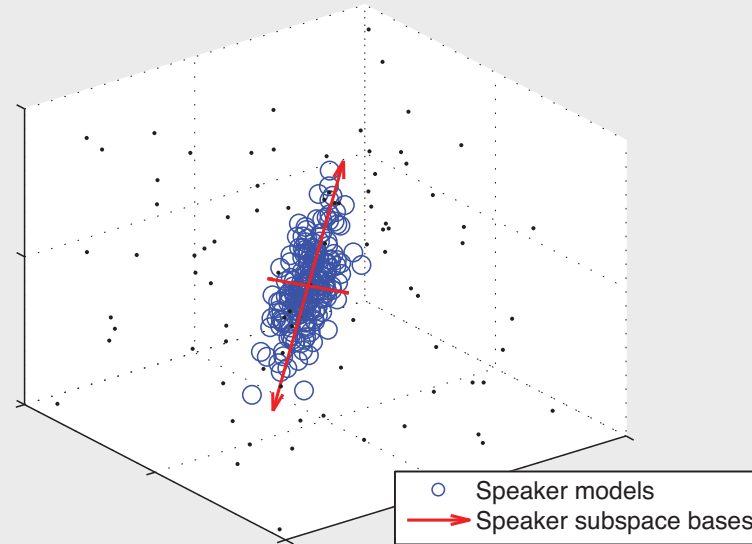  - exploit sequence constraints, speaker differences



- separation relies on detailed speaker model

# Eigenvoices

*Kuhn et al. '98, '00*
*Weiss & Ellis '10*

- Idea: Find **speaker model parameter space**

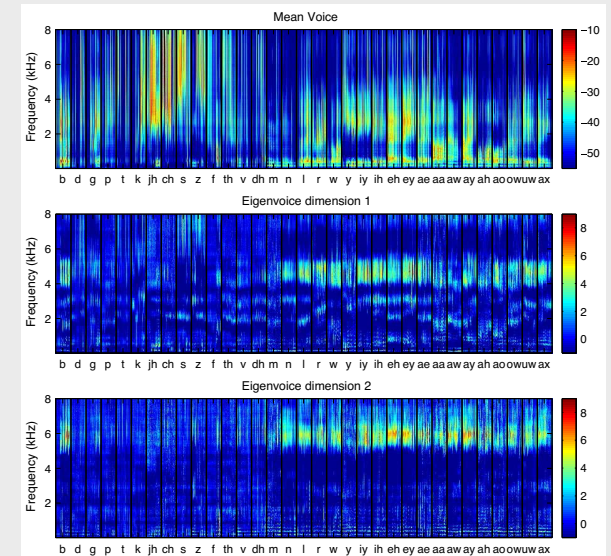  ○ generalize without losing detail?



○ Speaker models
→ Speaker subspace bases

- **Eigenvoice** model:

$$\boldsymbol{\mu} \; = \; \bar{\mu} \; + \; U \; \mathbf{w} \; + \; B \; \mathbf{h}$$

| adapted model | mean voice | eigenvoice bases | weights | channel bases | channel weights |

○ 89,600 dimensional space



Mean Voice

Eigenvoice dimension 1

Eigenvoice dimension 2

# Eigenvoice Speech Separation



Find Viterbi path

$$\boldsymbol{\mu}_1 = U\mathbf{w}_1 + \bar{\boldsymbol{\mu}}$$

$$\boldsymbol{\mu}_2 = U\mathbf{w}_2 + \bar{\boldsymbol{\mu}}$$

$\boldsymbol{y}(t)$

$\hat{\boldsymbol{x}}_1(t)$

$\hat{\boldsymbol{x}}_2(t)$

Update model parameters using EM algorithm from Kuhn et al., (2000)

Estimate source signals

# Eigenvoice Speech Separation

- **Eigenvoices for Speech Separation task**
  - speaker adapted (SA) performs midway between speaker-dependent (SD) & speaker-indep (SI)

# Spatial + Model Separation

- MESSL + Eigenvoice "priors"

**Observations**



Mixture – IPD
$\phi(\omega, t)$

Mixture – ILD
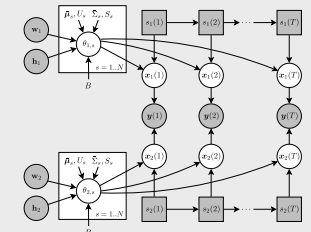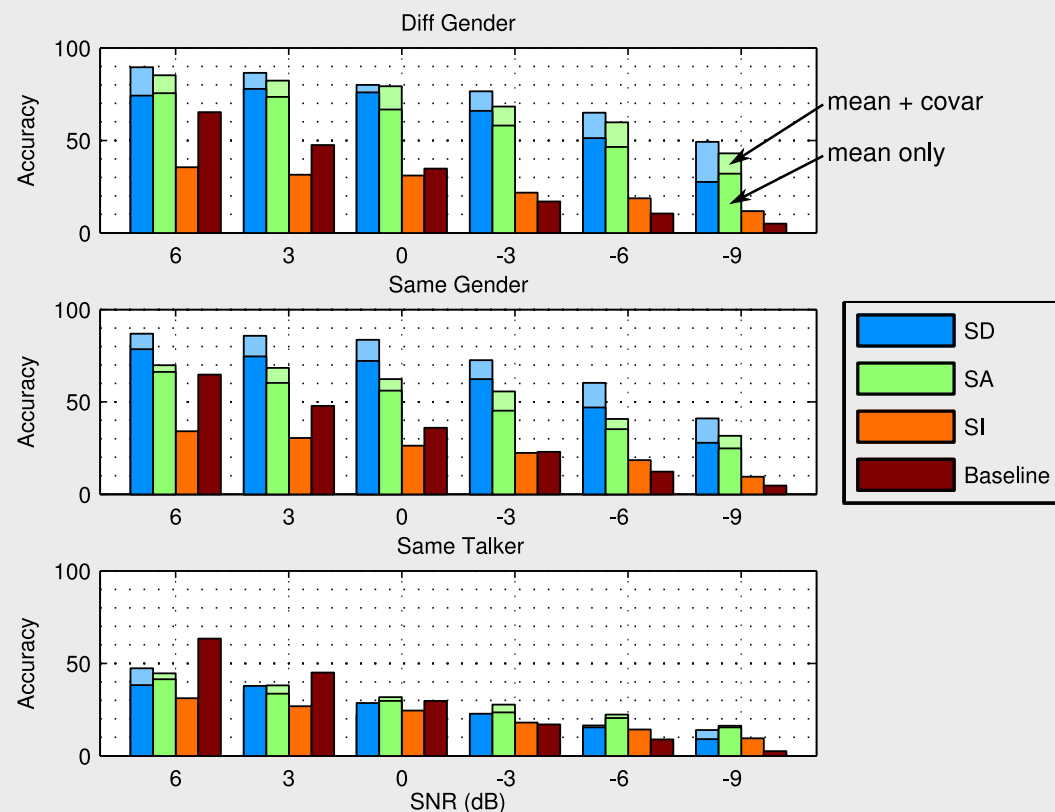$\alpha(\omega, t)$

Mixture – left channel
$\hat{L}(\omega, t)$

Mixture – right channel
$\hat{R}(\omega, t)$

~

**Parameters**

Per–source ITD
$\psi_{i\tau}$

Per–source ILD
$\mu_i, \eta_i$

~

Source prior (SP) means

SP covars
Mixture component

+

SP channel response – source 1
$h_1^\ell, h_1^r$

SP channel response – source 2
$h_2^\ell, h_2^r$

**Posteriors**

Each point in spectrogram is explained by a source, delay, and mixture component

**E-step**
Use parameters to compute posteriors of hidden variables

**M-step**
Use posteriors to update parameters

Source 1 mask

Source 2 mask

Separate sources by multiplying mixture by different masks

# Summary

- **Speech in the Wild**
  ... real, challenging problem
  ... applications in communications, lifelogs ...

- **Speech Separation**
  ... by generic properties (location, pitch)
  ... via statistical models

- **Recognition and Enhancement**
  ... separate-then-X, or integrated solution?

# References

- John Hershey, Steve Rennie, Pedr Olsen, Trausti Kristjansson, "Super-human multi-talker speech recognition: A graphical modeling approach," *Computer Speech & Lang.* 24 (1), 45-66, 2010.

- Jon Barker, Martin Cooke, Dan Ellis, "Decoding Speech in the Presence of Other Sources," *Speech Communication* 45(1): 5-25, 2005.

- R. Kuhn, J. Junqua, P. Nguyen, N. Niedzielski, "Rapid speaker adaptation in eigenvoice space," . *IEEE Tr. Speech & Audio Proc.* 8(6): 695–707, Nov 2000.

- Byung-Suk Lee & Dan Ellis, "Noise-robust pitch tracking by trained channel selection," submitted to *ICASSP*, 2012.

- Michael Mandel, Ron Weiss, Dan Ellis, "Model-Based Expectation-Maximization Source Separation and Localization," *IEEE Tr. Audio, Speech, Lang. Proc.* 18(2): 382-394, Feb 2010.

- A. Varga and R. Moore, "Hidden markov model decomposition of speech and noise," *ICASSP-90*, 845–848, 1990.

- Ron Weiss & Dan Ellis, "Speech separation using speaker-adapted Eigenvoice speech models," *Computer Speech & Lang.* 24(1): 16-29, 2010.

- Ron Weiss, Michael Mandel, Dan Ellis, "Combining localization cues and source model constraints for binaural source separation," *Speech Communication* 53(5): 606-621, May 2011.

- Mingyang Wu, DeLiang Wang, Guy Brown, "A multipitch tracking algorithm for noisy speech," *IEEE Tr. Speech & Audio Proc.* 11(3): 229–241, May 2003.