# Using Speech Models for Separation

## Dan Ellis

*Comprising the work of Michael Mandel and Ron Weiss*
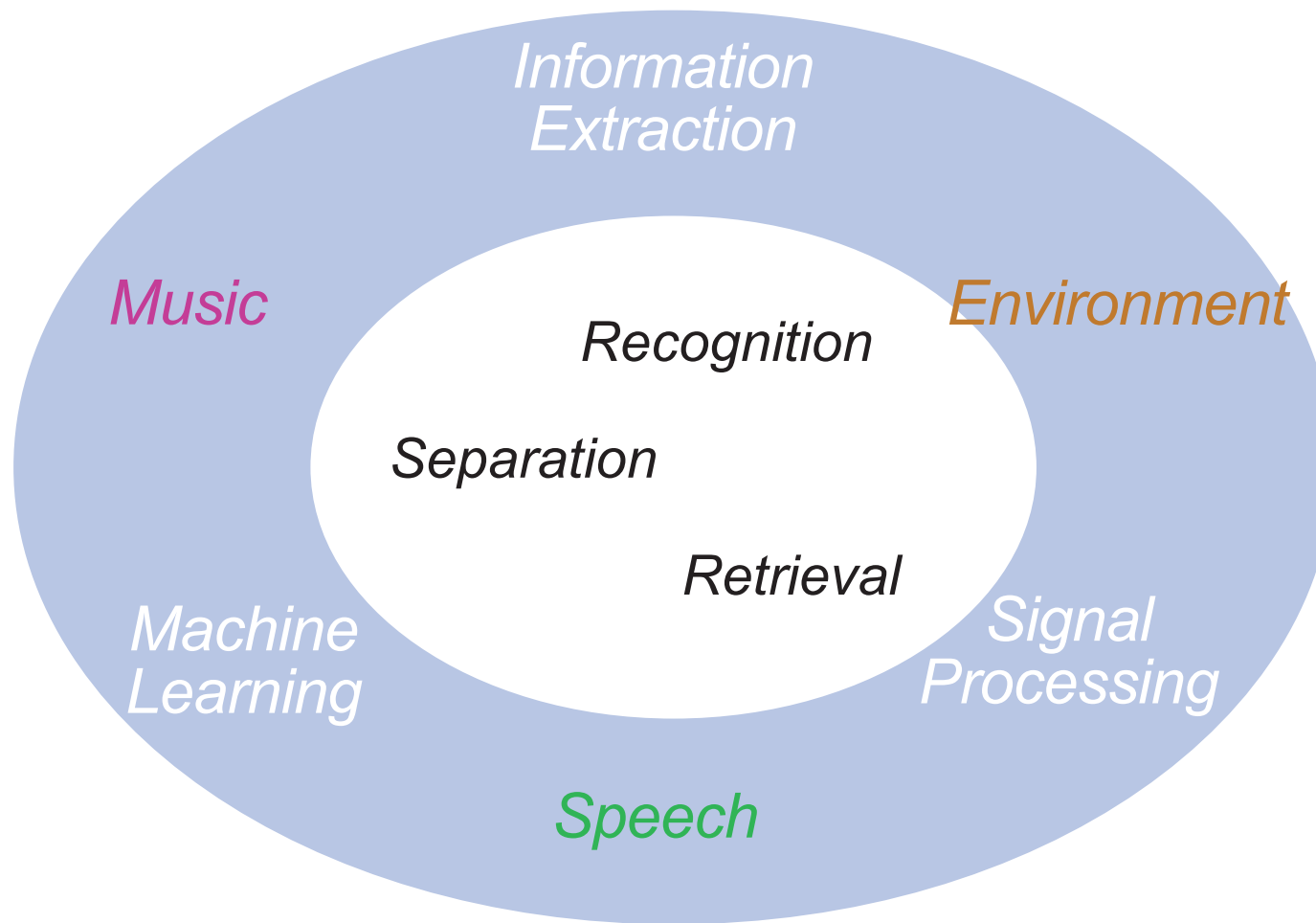Laboratory for Recognition and Organization of Speech and Audio
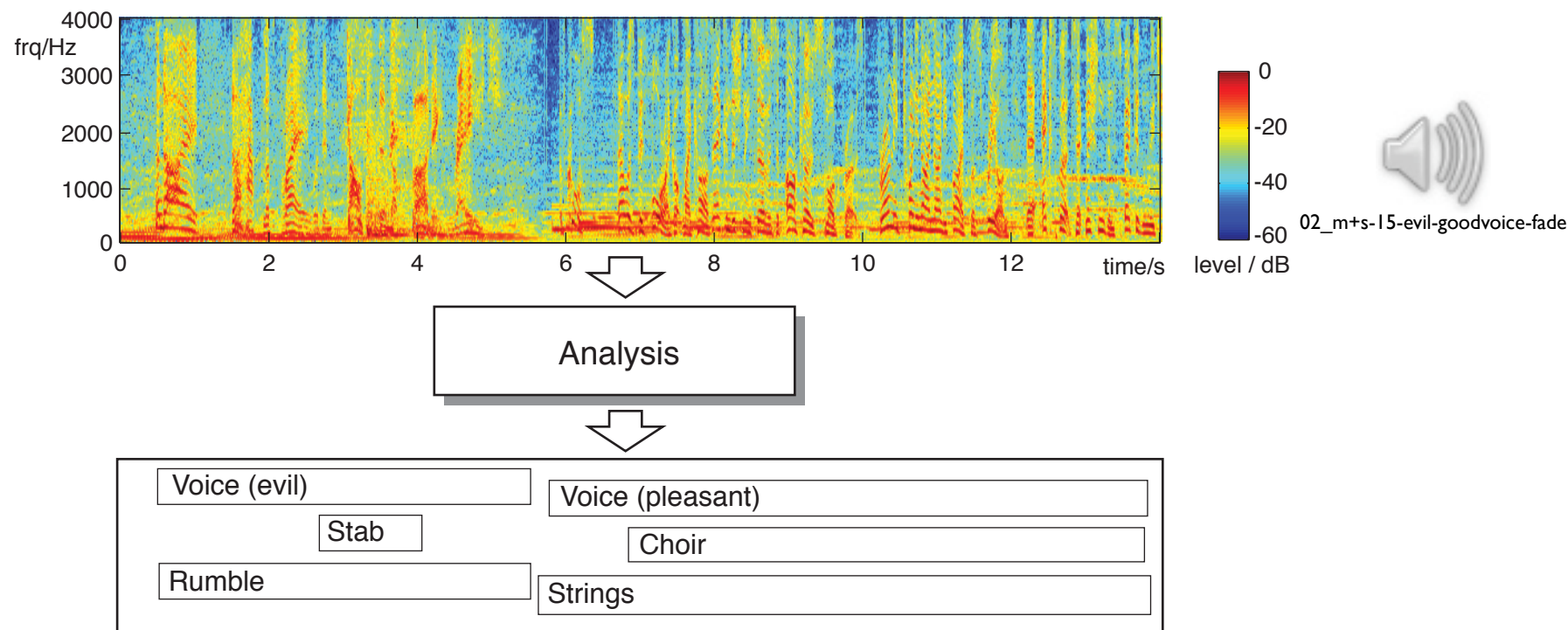Dept. Electrical Eng., Columbia Univ., NY USA

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

# LabROSA Overview



Information Extraction

Music

Environment

Recognition

Separation

Retrieval

Machine Learning

Signal Processing

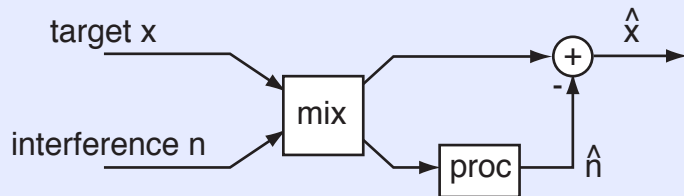Speech

# 1. Source Models and Scene Analysis



- **Sounds rarely occur in isolation**
  - .. so analyzing mixtures ("scenes") is a problem
  - .. for humans and machines

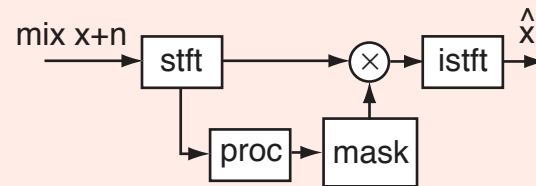# Approaches to Separation

## ICA

- Multi-channel
- Fixed filtering
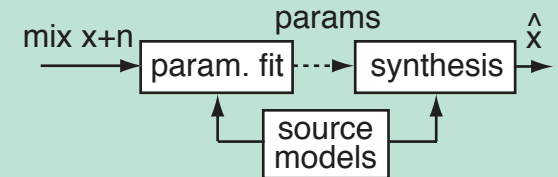- Perfect separation – maybe!



## CASA

- Single-channel
- Time-var. filter
- Approximate separation



## Model-based

- Any domain
- Param. search
- Synthetic output?
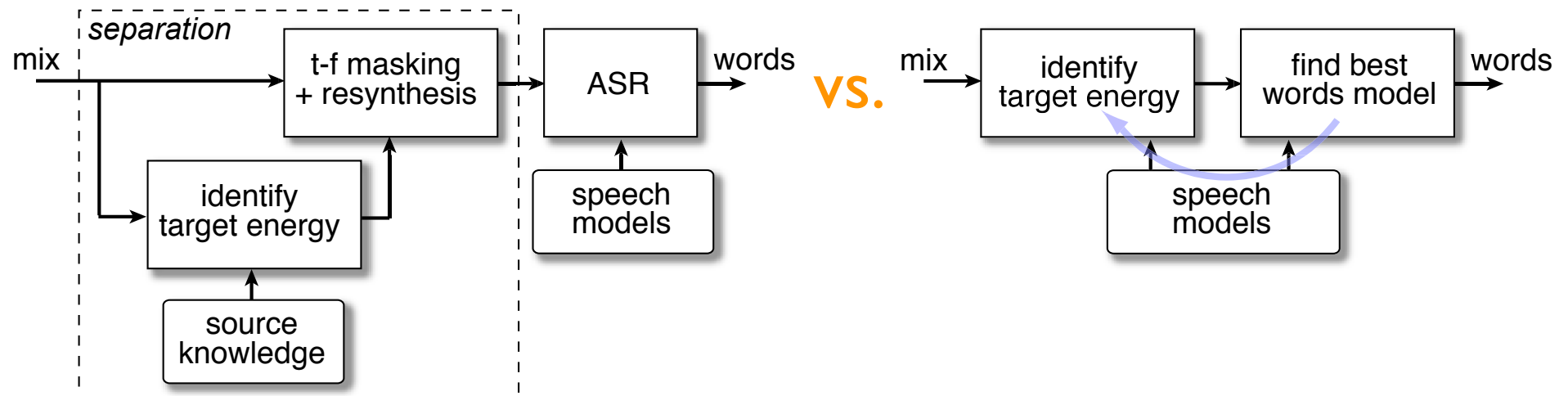
# Separation vs. Inference

- **Ideal** separation is rarely possible
  - many situations where overlaps cannot be removed

- Overlaps → Ambiguity
  - scene analysis = find "most reasonable" explanation

- Ambiguity can be expressed probabilistically
  - i.e. posteriors of sources $\{S_i\}$ given observations $X$:

$$P(\{S_i\}\mid X) \propto \underbrace{P(X \mid \{S_i\})}_{\text{combination physics}} \prod_i \underbrace{P(S_i \mid M_i)}_{\text{source models}}$$

  - search over all source signal sets $\{S_i\}$ ??

- Better source models → better inference

# 2. Speech Separation Using Models

- **Cooke & Lee's Speech Separation Challenge**
  - pairs of short, grammatically-constrained utterances:
    `<command:4><color:4><preposition:4><letter:25><number:10><adverb:4>`
    e.g. "bin white by R 8 again"
  - task: report letter + number for "white"
  - (special session at Interspeech '06)
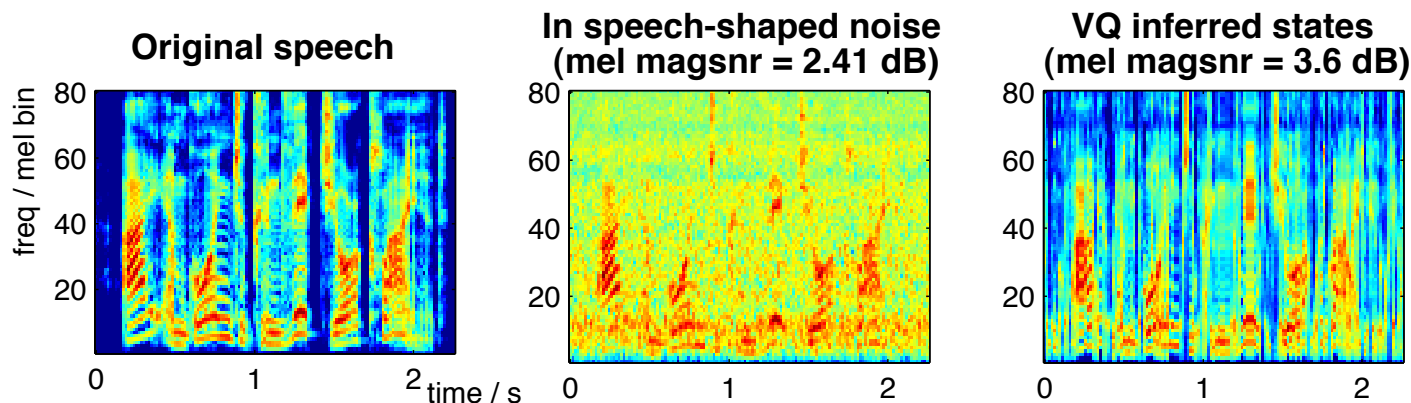
- **Separation or Description?**

# Codebook Models

- Given models for sources,
  find "best" (most likely) states for spectra:

$$p(\mathbf{x}|i_1, i_2) = \mathcal{N}(\mathbf{x}; \mathbf{c}_{i1} + \mathbf{c}_{i2}, \Sigma)$$ *combination model*

$$\{i_1(t), i_2(t)\} = argmax_{i_1, i_2} p(\mathbf{x}(t)|i_1, i_2)$$ *inference of source state*

  ○ can include sequential constraints...

- E.g. stationary noise:



**Original speech**

**In speech-shaped noise (mel magsnr = 2.41 dB)**

**VQ inferred states (mel magsnr = 3.6 dB)**

# Speech Recognition Models

*Varga & Moore '90*

- **Speech recognizers contain speech models**
  - ASR is just $\mathrm{argmax}\ \mathrm{P}(W \mid X)$
- **Recognize mixtures with Factorial HMM**
  - i.e. two state sequences, one model for each voice
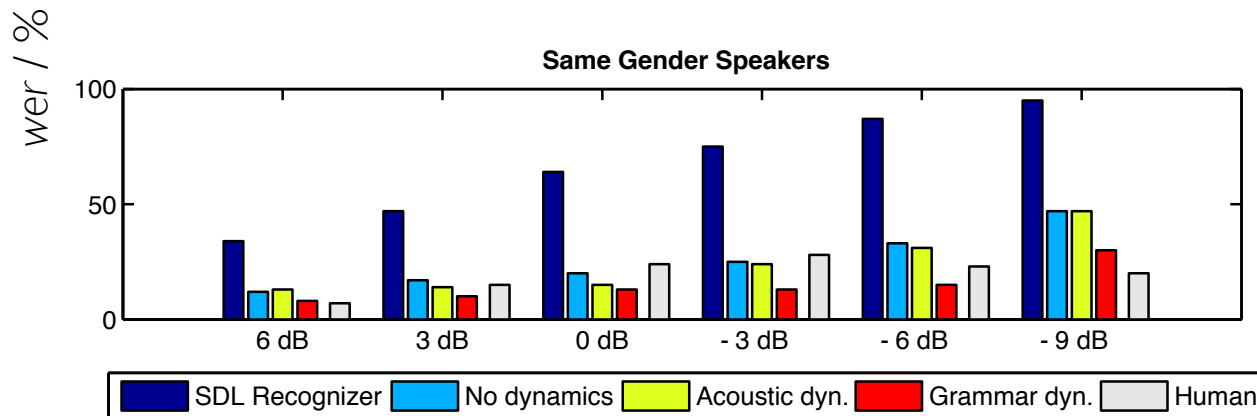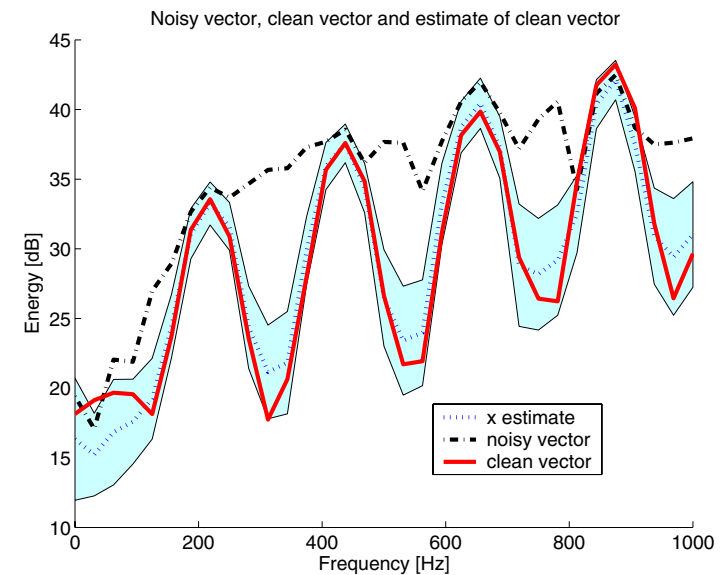  - exploit sequence constraints, speaker differences

# Speech Factorial Separation

- IBM's 2006 Iroquois speech separation system
  Key features:
  - detailed state combinations
  - large speech recognizer
  - exploits grammar constraints
  - 34 per-speaker models

- "Superhuman" performance
  - ... in some conditions



Noisy vector, clean vector and estimate of clean vector

Legend:
- x estimate
- noisy vector
- clean vector



Same Gender Speakers

Legend: SDL Recognizer, No dynamics, Acoustic dyn., Grammar dyn., Human

# Adapting Source Models

- Power of model-based separation depends on detail of model
- Speech separation relies on prior knowledge of every speaker?



- Can this be practical?

# Eigenvoices

- Idea: Find
**model parameter space**

  ○ generalize without
  losing detail?



- Speaker models
- Speaker subspace bases

- **Eigenvoice** model:

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

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

# Eigenvoice Bases

- **Mean model**
  - 280 states x 320 bins = 89,600 dimensions

- **Eigencomponents shift formants/ coloration**
  - additional components for channel

# Speaker-Adapted Separation

- Factorial HMM analysis
  with tuning of source model parameters
  = eigenvoice speaker adaptation
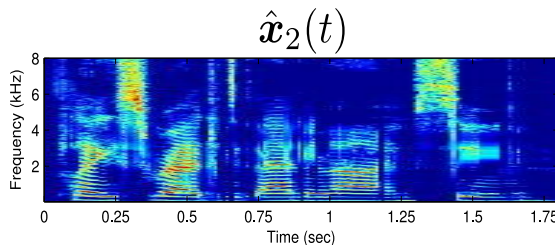
# Speaker-Adapted Separation

Find Viterbi path

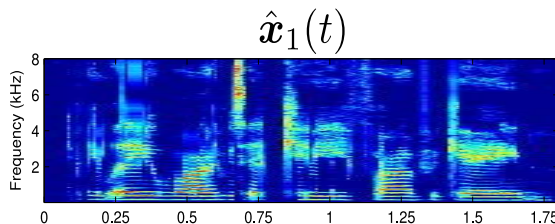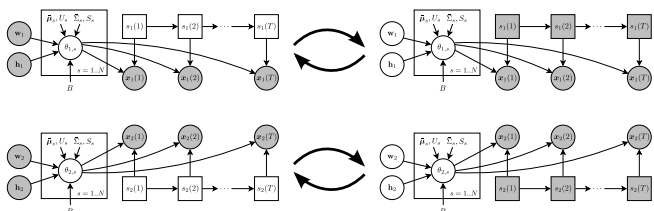$\boldsymbol{\mu}_1 = U\mathbf{w}_1 + \bar{\boldsymbol{\mu}}$

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

$\boldsymbol{y}(t)$

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

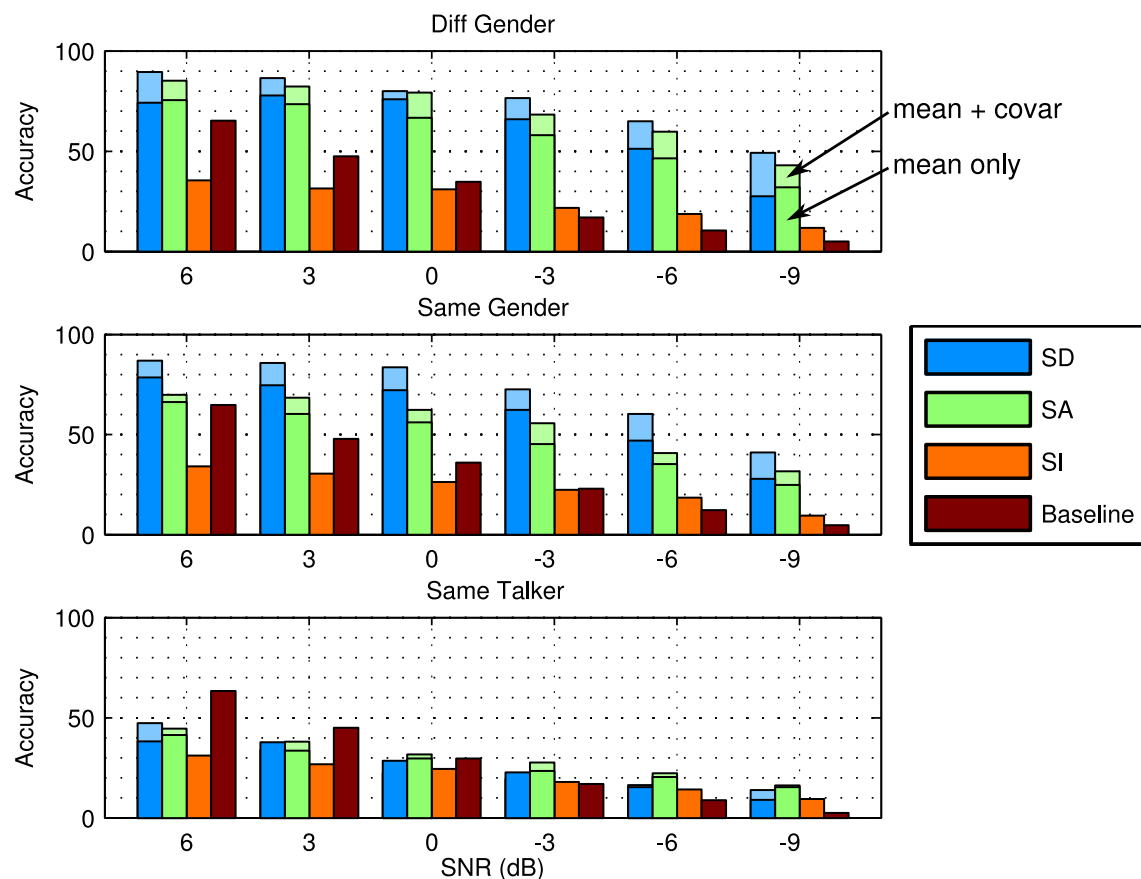$\hat{\boldsymbol{x}}_1(t)$

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

Estimate source signals

# Speaker-Adapted Separation

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

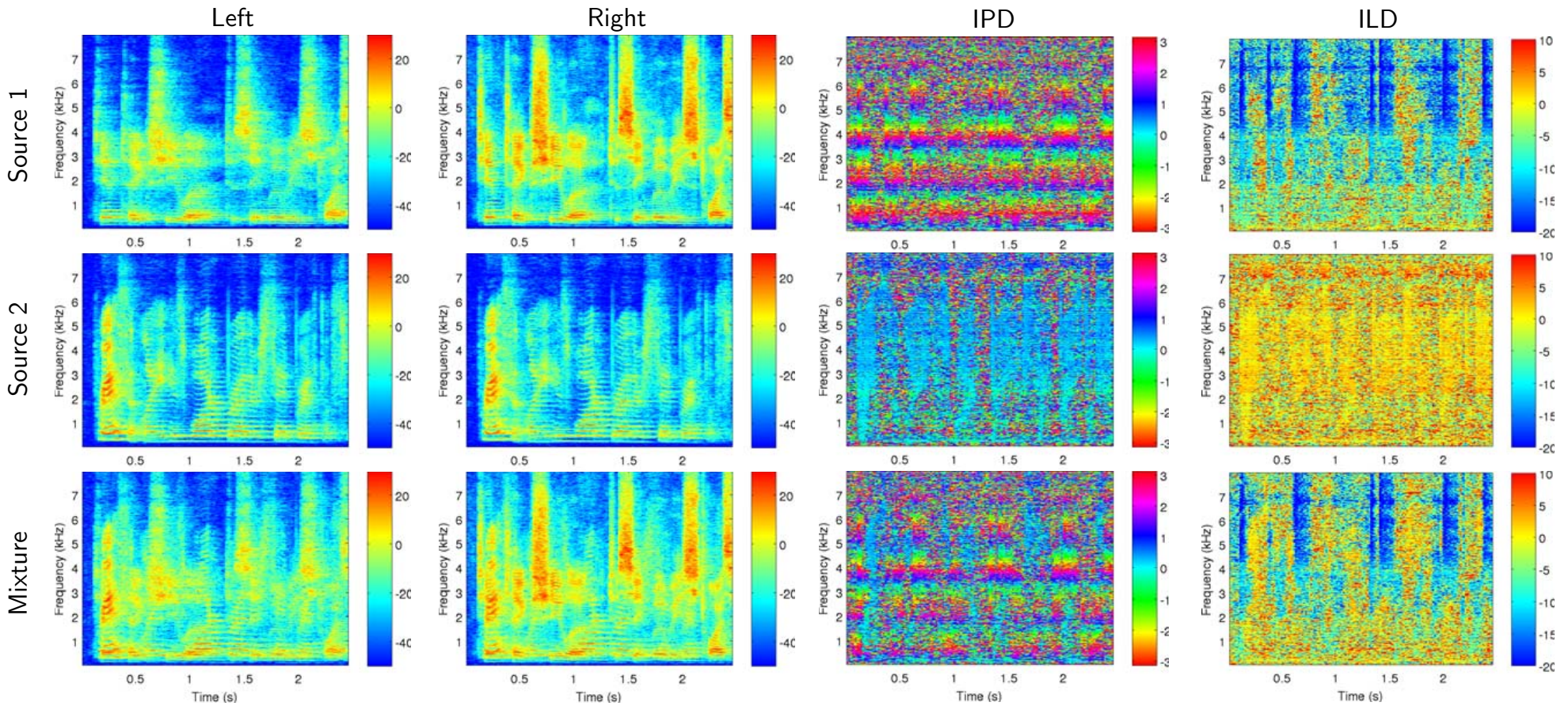- 2 or 3 sources
  in reverberation

  o assume just 2 'ears'



- Model interaural spectrum of each source
  as stationary level and time differences:

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

# ILD and IPD

- Sources at 0° and 75° in reverb
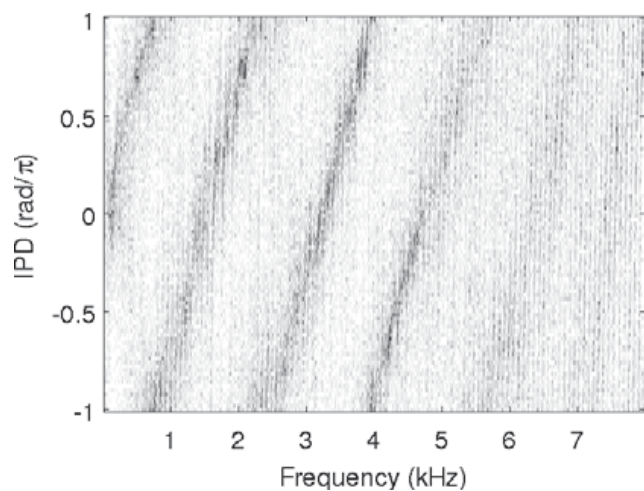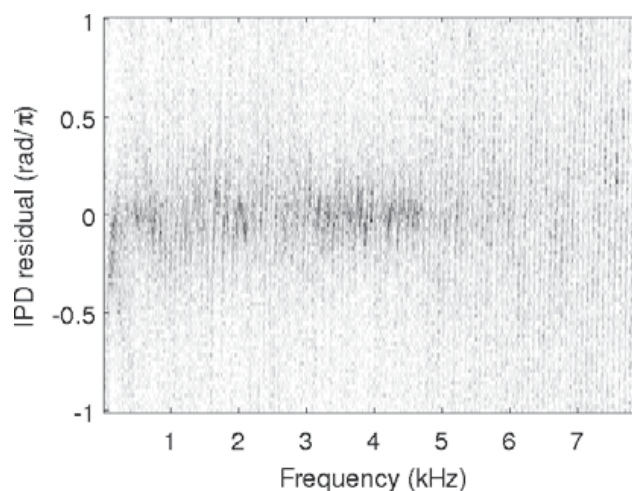
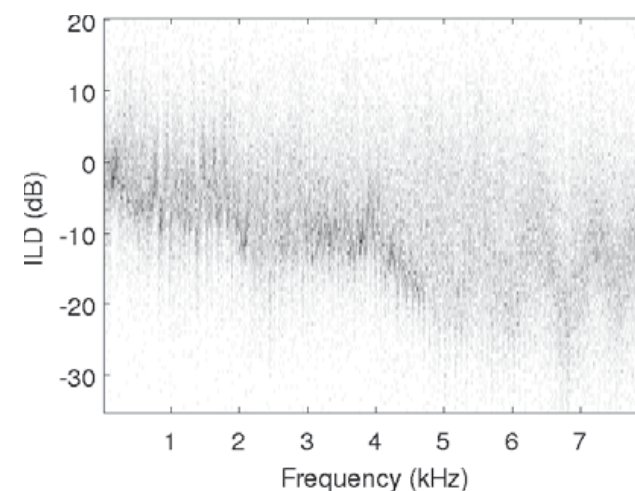# IPD, ILD Distributions

- Source at 75° in reverberation

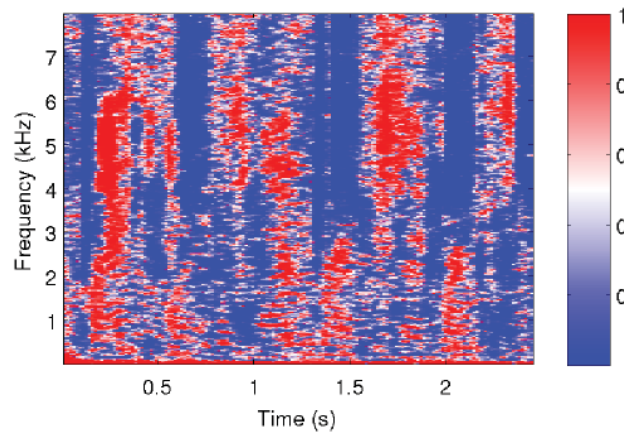IPD             IPD residual             ILD



- IPD residual offsets phase by constant $\omega\tau$
- IPD can be fit by single Gaussian
- ILD needs frequency-dependence

# Model-Based EM Source Separation and Localization (MESSL)

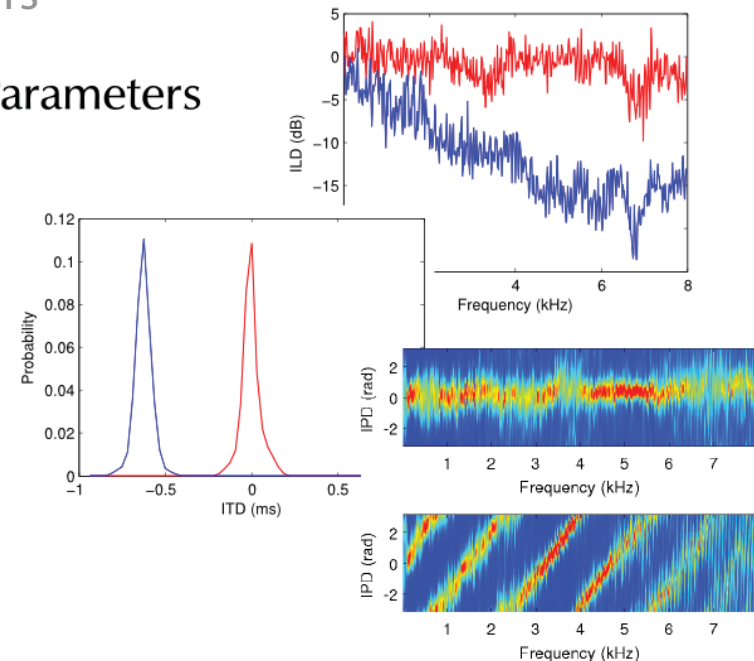*Mandel & Ellis '09*

Re-estimate
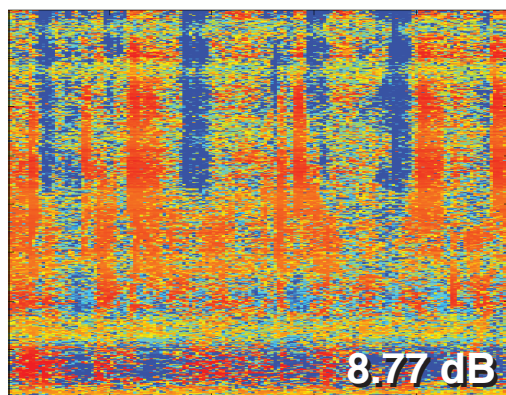source parameters
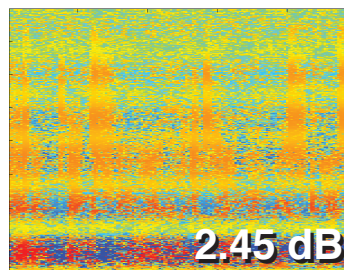


Assign spectrogram points
to sources

- can model more sources than sensors
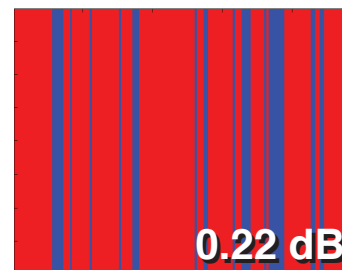- flexible initialization

# MESSL Results

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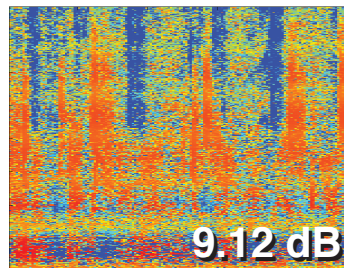  - tradeoff between constraints & noisiness



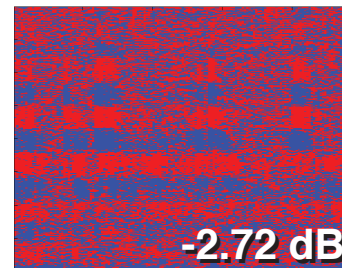EM+ILD — 8.77 dB

EM-ILD (only IPD) — 2.45 dB

EM+1ILD (tied means) — 9.12 dB

PHAT-histogram — 0.22 dB

DUET — -2.72 dB

Ground Truth — 12.35 dB

# MESSL Results

- ## Signal-to-Distortion Ratio (SDR)



Anechoic 2 talkers — Reverb 2 talkers — Reverb 3 talkers

Legend: DP–Wiener, Wiener, MESSL–G, MESSL–$\Omega\Omega$, TRINICON, Mouba, Sawada, DUET, Random
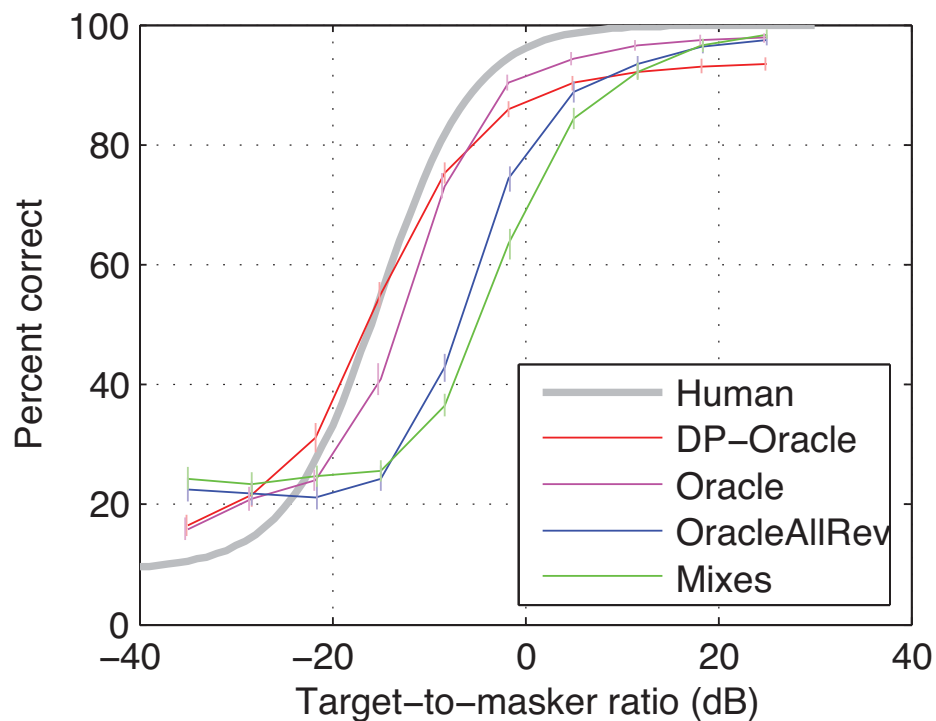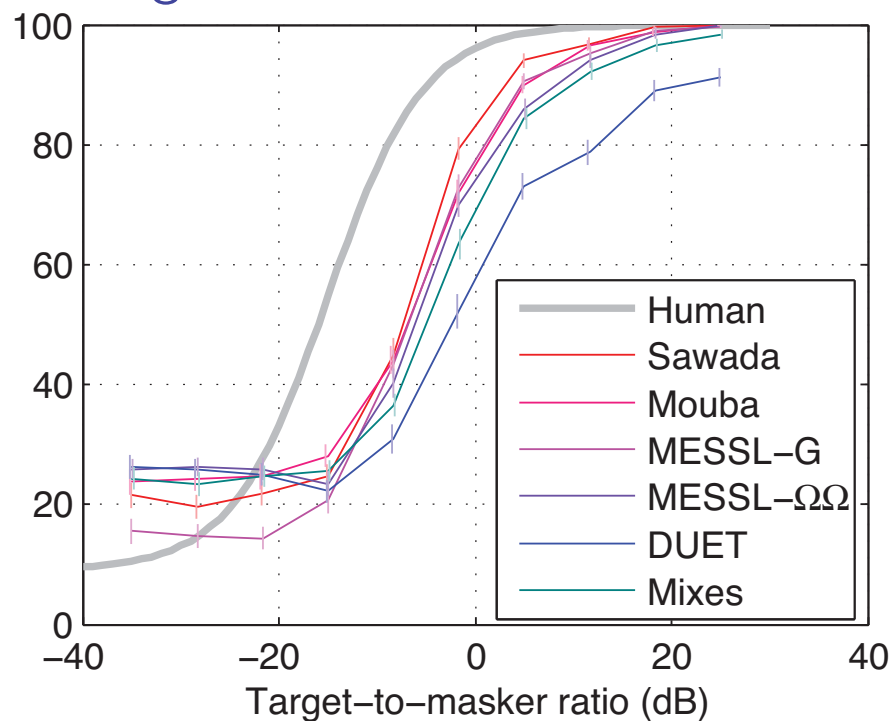
# MESSL Results

- Speech recognizer (Digits)



Ground truth masks

Algorithmic masks

# 4. Combining Spatial + Speech Models

- Interaural parameters give

$$ILD_i(\omega),\ ITD_i,\ \Pr\big(X(t,\omega) = S_i(t,\omega)\big)$$

- Speech source model can give

$$\Pr\big(S_i(t,\omega) \text{ is speech signal}\big)$$

- Can combine into one big EM framework...

**E-step**
$$p(u|\Theta^{(n)}) = p(x,u|\Theta^{(n)})/p(x|\Theta^{(n)})$$

*u is: Pr(cell from source i)*
*phoneme sequence*

**M-step**
$$\Theta^{(n+1)} = \underset{\Theta}{\operatorname{argmax}}\ E_{p(u|\Theta^{(n)})}\, p(x,u|\Theta)$$

*Θ is: interaural params*
*speaker params*

# MESSL-SP (Source Prior)

**Observations**

Mixture – IPD

$\phi(\omega, t)$

Mixture – ILD

$\alpha(\omega, t)$

Mixture – left channel

$\hat{L}(\omega, t)$

Mixture – right channel

$\hat{R}(\omega, t)$

Time (sec)

~

~

~

**Parameters**

Per–source ITD

$\psi_{i\tau}$

Time (samples)

Per–source ILD

$\mu_i, \eta_i$

Frequency (kHz)

Source prior (SP) means

SP covars

Mixture component

+

SP channel response – source 1

$h_1^\ell, h_1^r$

SP channel response – source 2

$h_2^\ell, h_2^r$

Frequency (kHz)

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

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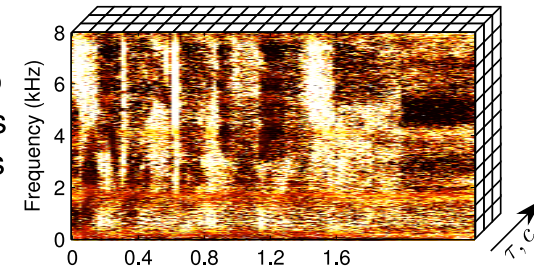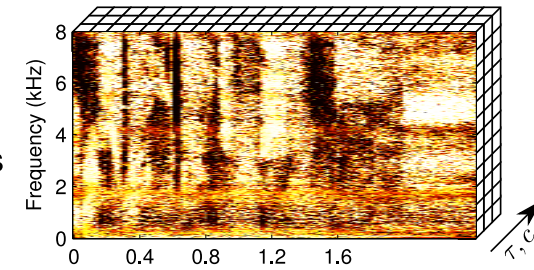Use posteriors to
update parameters

**Posteriors**

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

Source 1 mask

Source 2 mask

Time (sec)

Separate sources by
multiplying mixture
by different masks

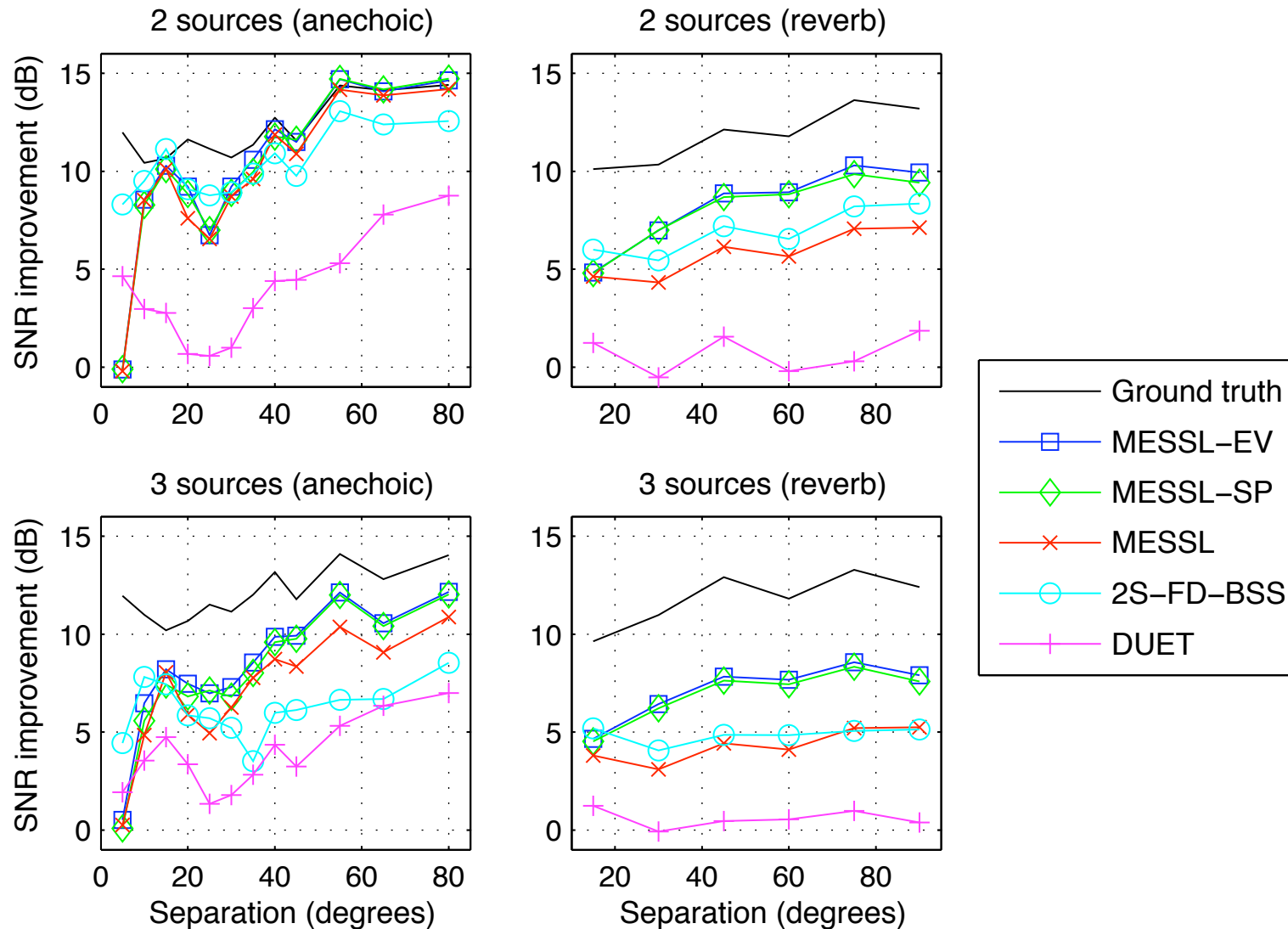# MESSL-SP Results

- Source models function as priors
- Interaural parameter spatial separation
  - source model prior improves spatial estimate



Ground truth (12.04 dB) — DUET (3.84 dB) — 2D–FD–BSS (5.41 dB)
MESSL (5.66 dB) — MESSL–SP (10.01 dB) — MESSL–EV (10.37 dB)

# MESSL-SP Results

- SNR improvement vs. source angle separation

# Future Work

- **Better** parametric speaker models
  - limitations of eigenvoices
  - varying style

- **Understanding** reverb & ASR
  - early echoes
  - what spoils ASR?

- **Models of** other sources
  - eigeninstruments?

# Summary & Conclusions

- **Source models** provide the constraints to make scene analysis possible

- **Eigenvoices** (model subspace) can be used to provide detailed models that generalize

- Spatial parameters can identify more sources than models in reverb (MESSL)

- Can combine source + spatial models