

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

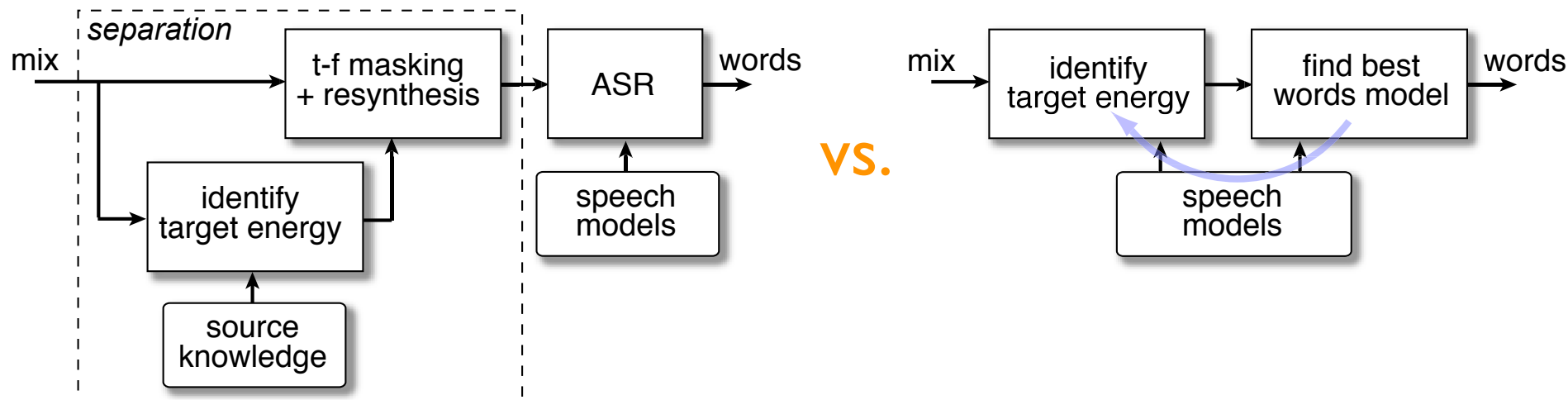
dpwe@ee.columbia.edu

<http://labrosa.ee.columbia.edu/>

1. Eigenvoice Speaker Models
2. Spatial Parameter Models in Reverb
3. Combining Source + Spatial

I. Speech Separation Using Models

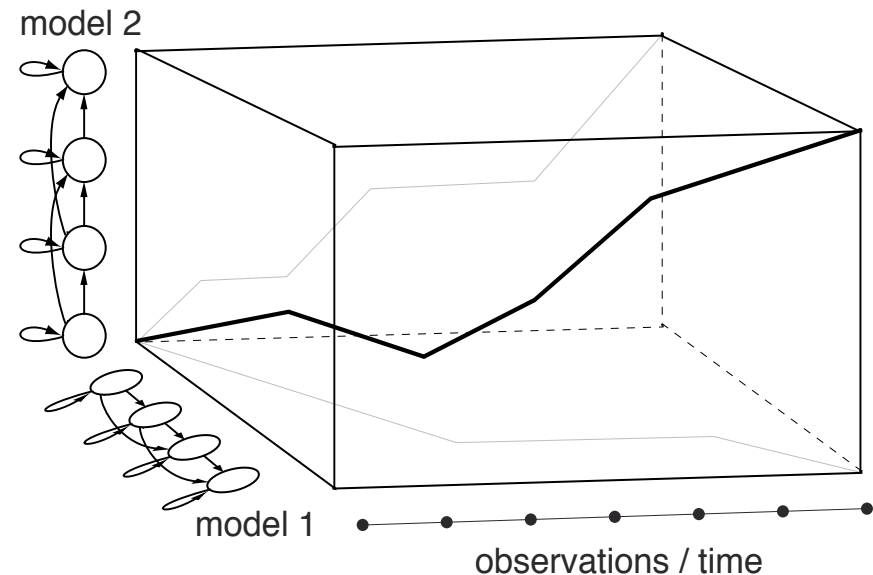
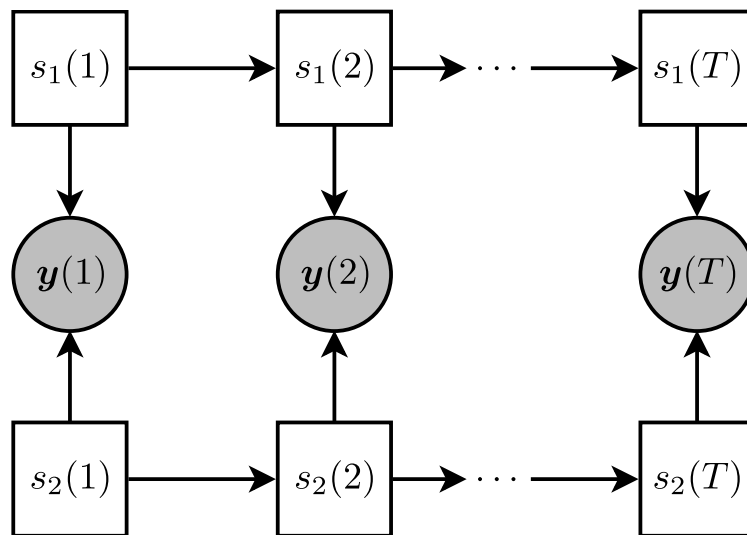
- **Cooke & Lee's Monaural Speech Separation Task**
 - pairs of short, grammatically-constrained utterances:
<command:4><color:4><preposition:4><letter:25><number:10><adverb:4>
e.g. "bin white by R 8 again"
 - task: report letter + number for "white"
- Separation depends on **source constraints**
 - the more the better - ASR model



Speech Mixture Recognition

Kristjansson, Hershey et al. '06

- Speech recognizers contain speech models
 - ASR is just $\operatorname{argmax} P(W | X)$
- Recognize mixtures with **Factorial HMM**
 - one model+state sequence for each voice
 - exploit sequence constraints, **speaker differences**



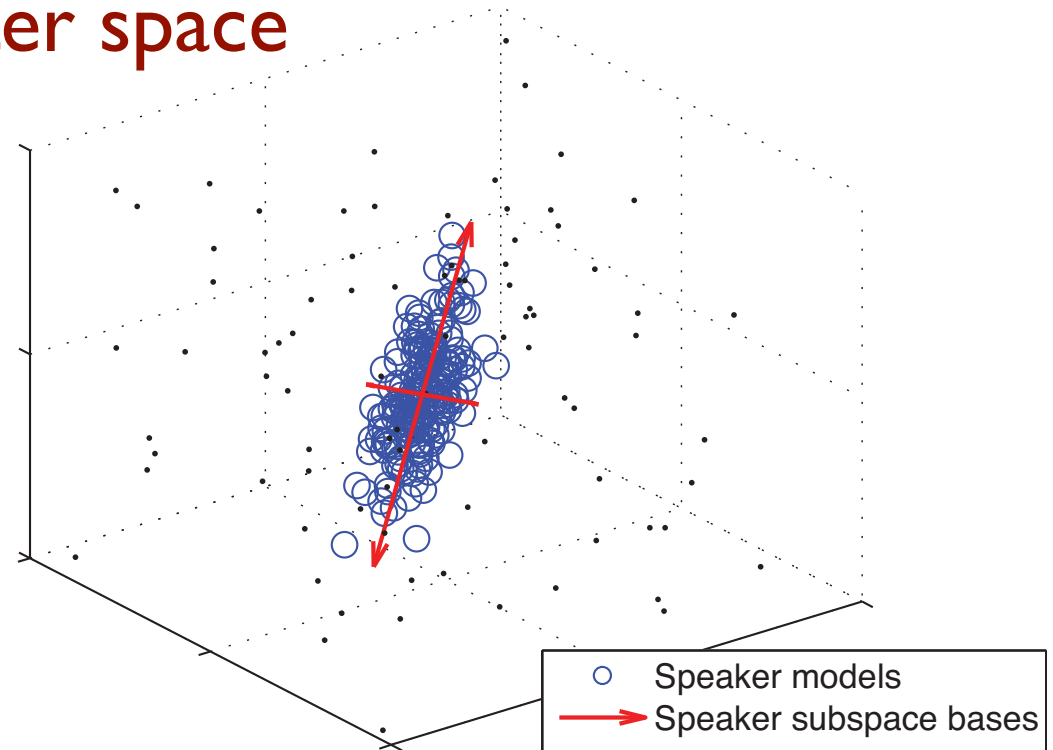
- separation relies on **detailed speaker model**

Eigenvoices

Kuhn et al. '98, '00
Weiss & Ellis '07, '08, '09

- Idea: Find speaker model parameter space

- generalize without losing detail?



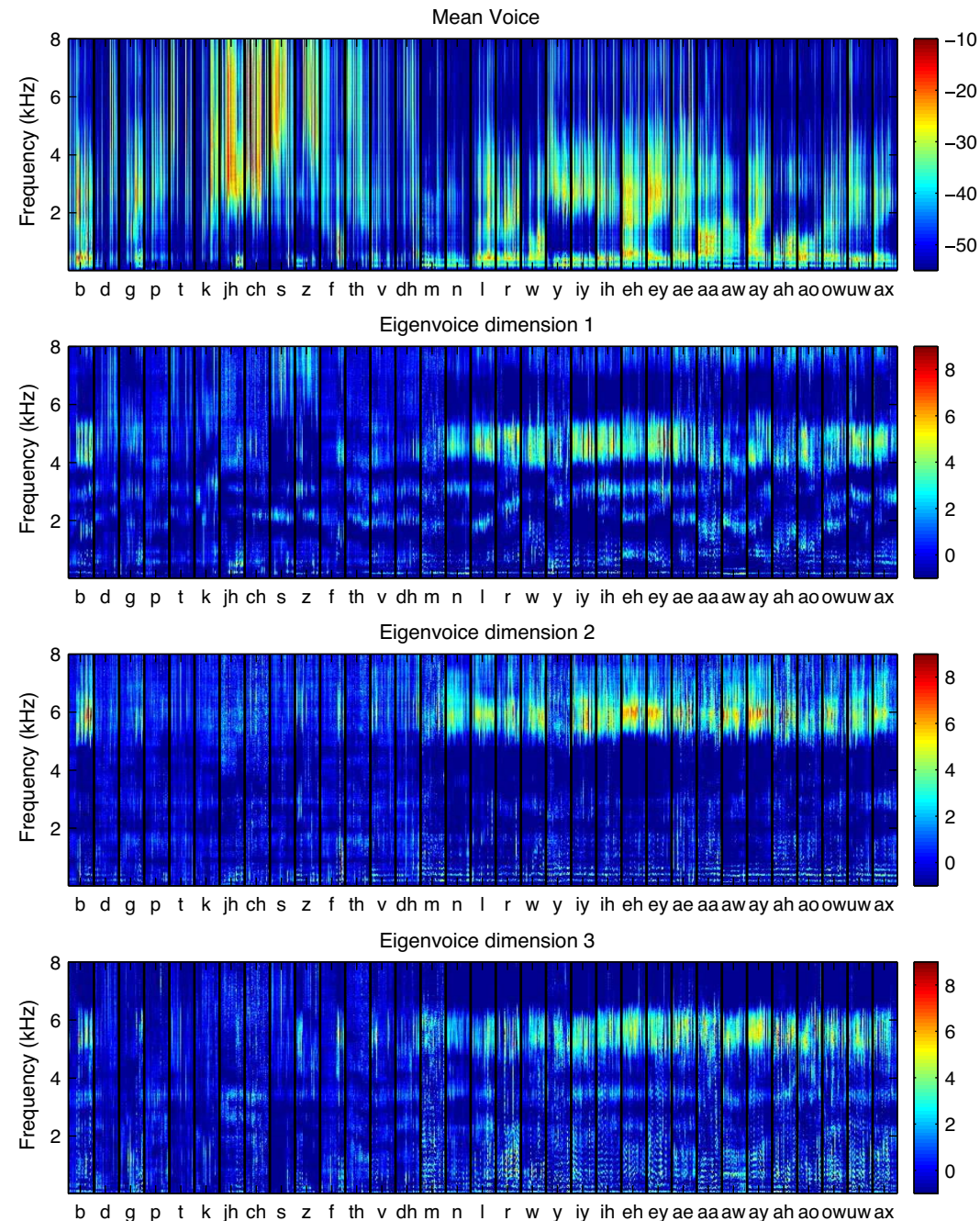
- Eigenvoice model:

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

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

Eigenvoice Bases

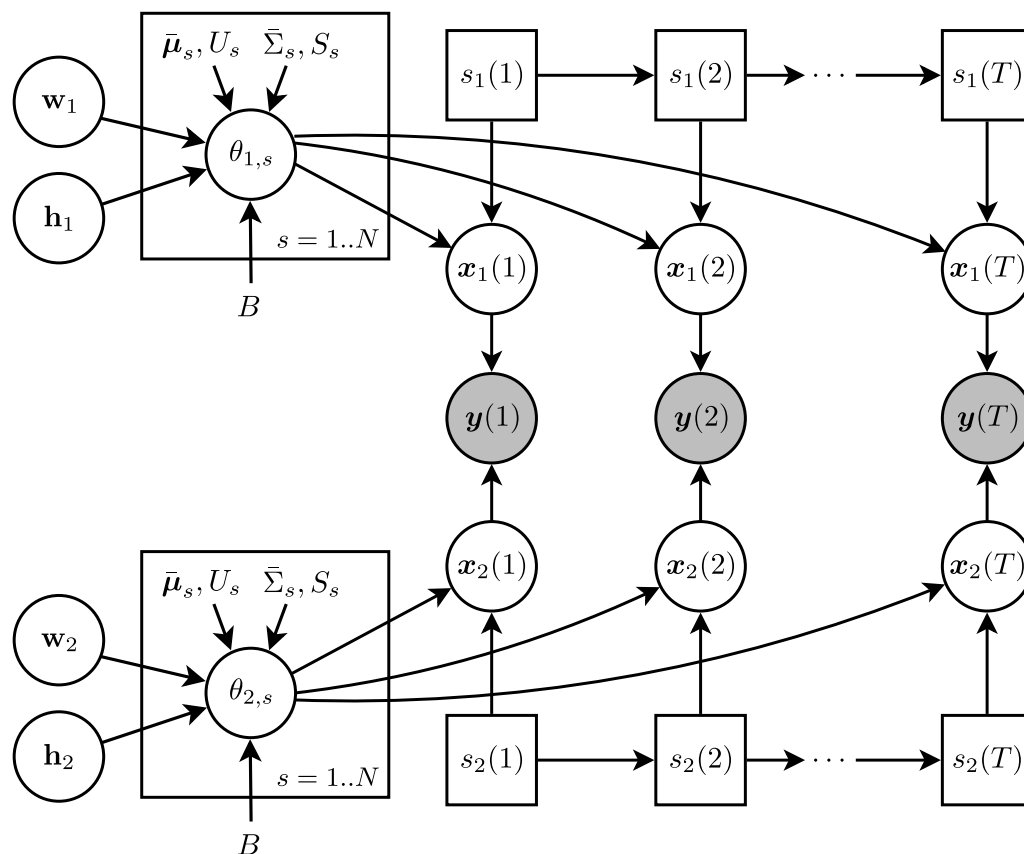
- Mean model
 - 280 states \times 320 bins
= 89,600 dimensions
- Eigencomponents shift formants/
coloration
 - additional components for acoustic channel



Speaker-Adapted Separation

Weiss & Ellis '08

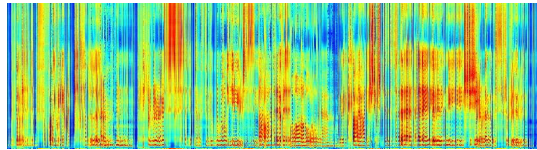
- Factorial HMM analysis
with **tuning** of source model parameters
= **eigenvoice speaker adaptation**



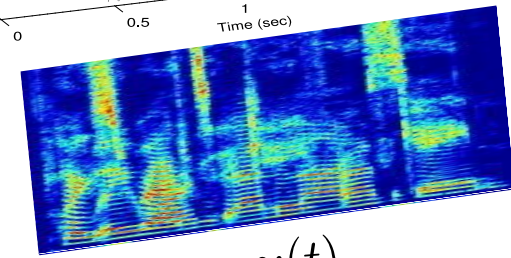
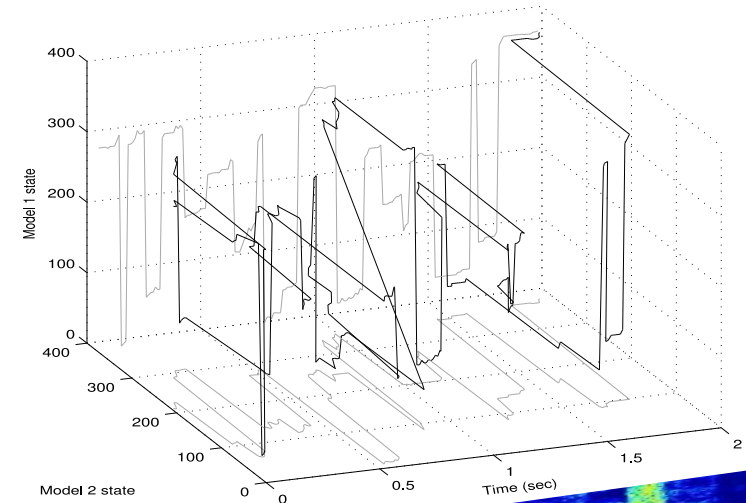
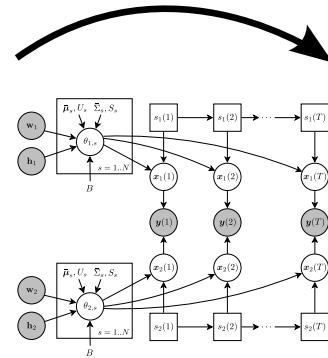
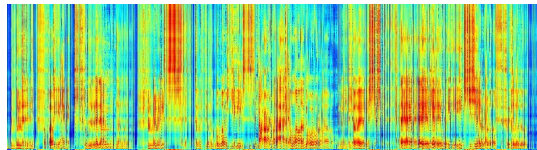
Speaker-Adapted Separation

Find Viterbi path

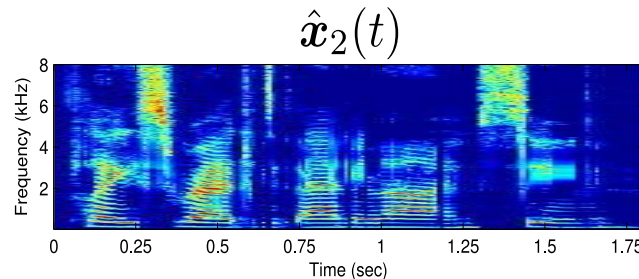
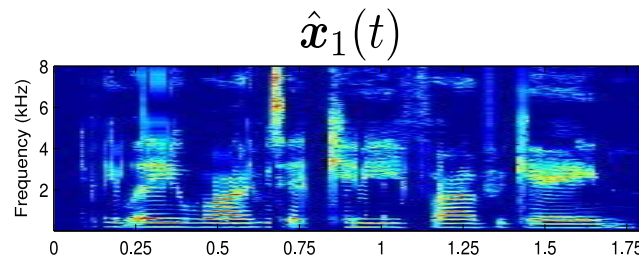
$$\mu_1 = U\mathbf{w}_1 + \bar{\mu}$$



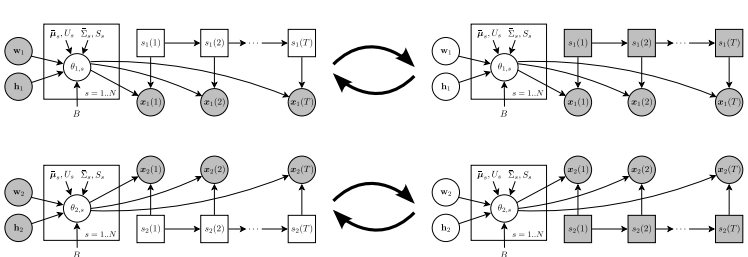
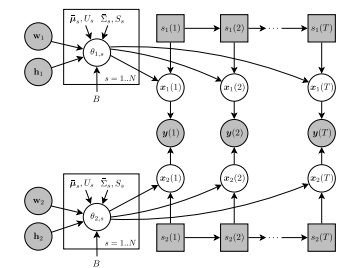
$$\mu_2 = U\mathbf{w}_2 + \bar{\mu}$$



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

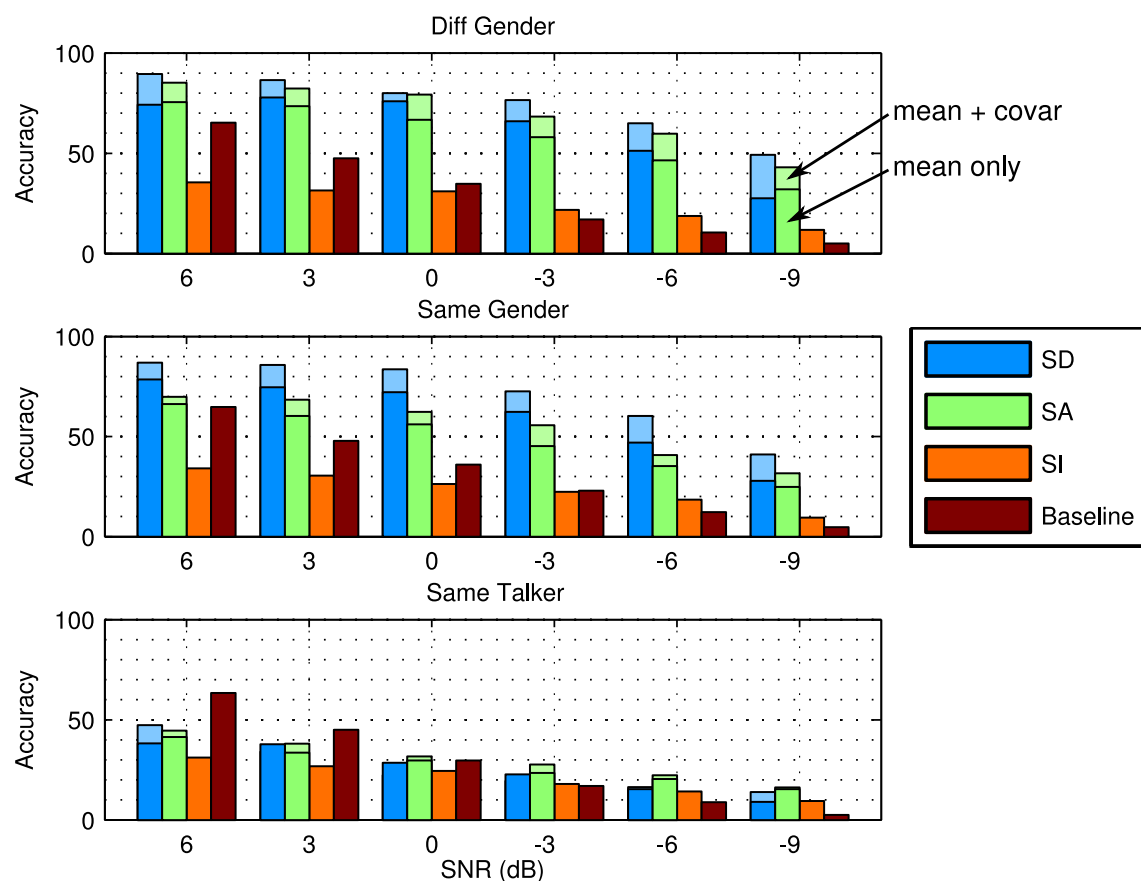


Estimate source signals



Speaker-Adapted Separation

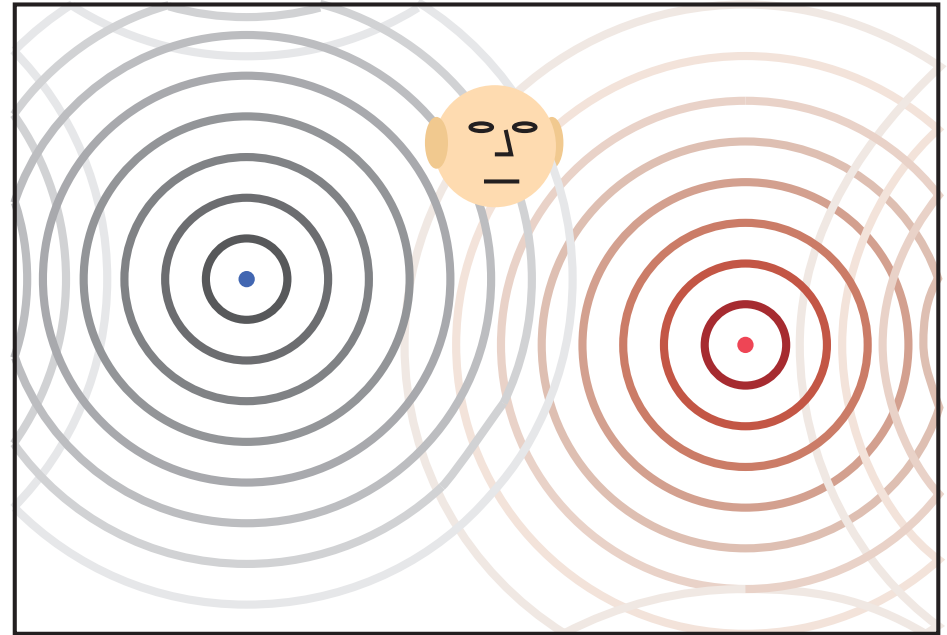
- Eigenvoices for Speech Separation task
 - speaker adapted (SA) performs midway between speaker-dependent (SD) & speaker-indep (SI)



2. Spatial Models & Reverb

Mandel & Ellis '07

- 2 or 3 sources in reverberation
 - assume just 2 'ears'



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

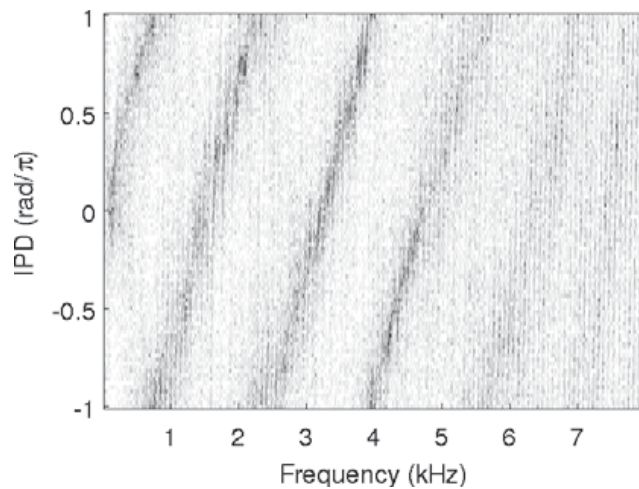
$$\frac{L(\omega, t)}{R(\omega, t)} = a(\omega) e^{j\omega\tau} N(\omega, t)$$

IPD, ILD Distributions

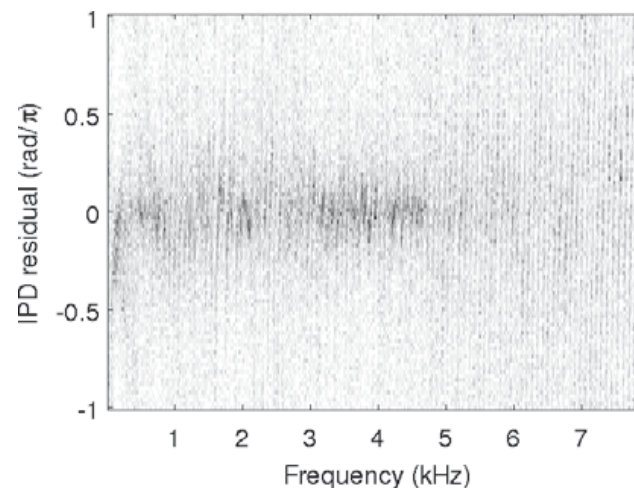
- Source at 75° in reverberation



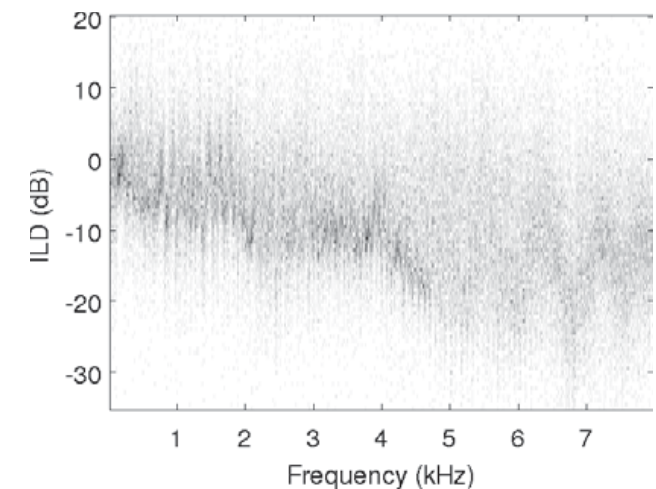
IPD



IPD residual



ILD



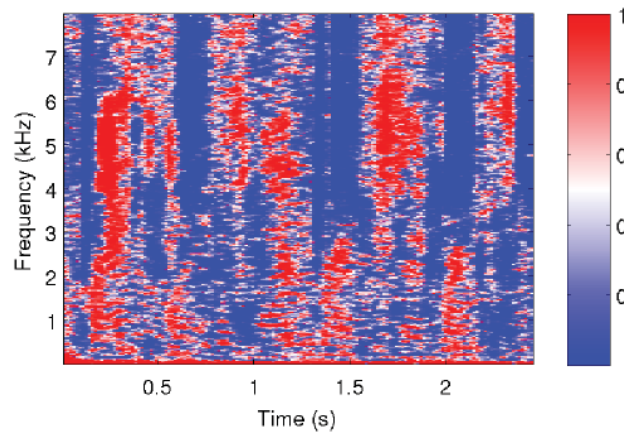
- IPD residual offsets phase by constant ωT
- IPD can be fit by single Gaussian
- ILD needs frequency-dependence

Model-Based EM Source Separation and Localization (MESSL)

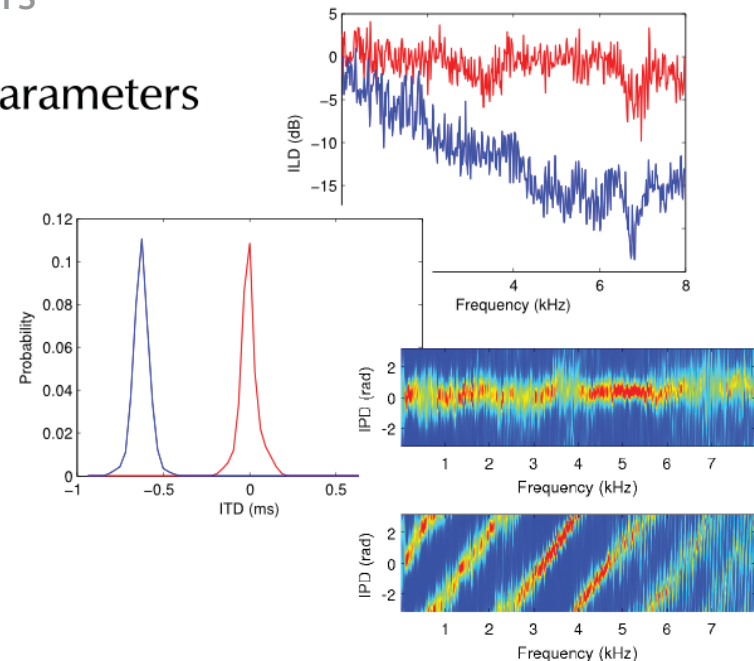
Mandel & Ellis '09

Re-estimate
source parameters

Masks



Parameters

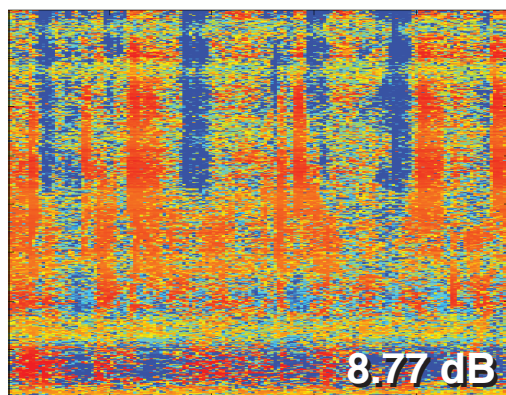


Assign spectrogram points
to sources

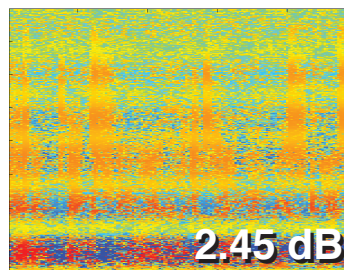
- can model more sources than sensors
- flexible initialization

MESSL Results

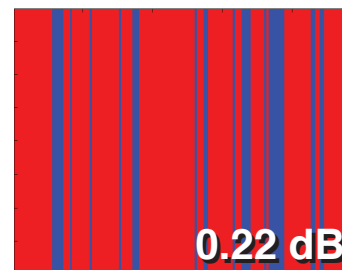
- **Modeling uncertainty** improves results
 - tradeoff between constraints & **noisiness**



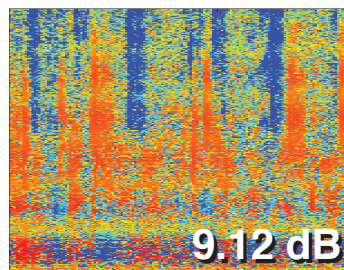
EM+ILD



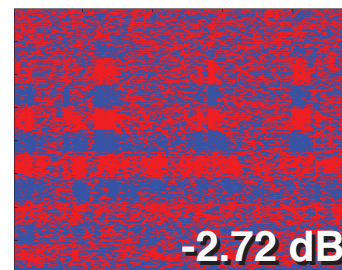
EM-ILD (only IPD)



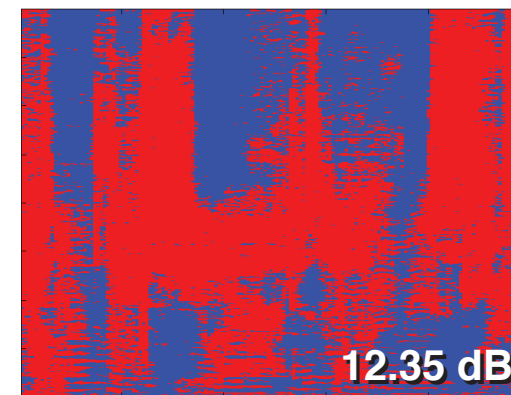
PHAT-histogram



EM+1ILD (tied means)



DUET

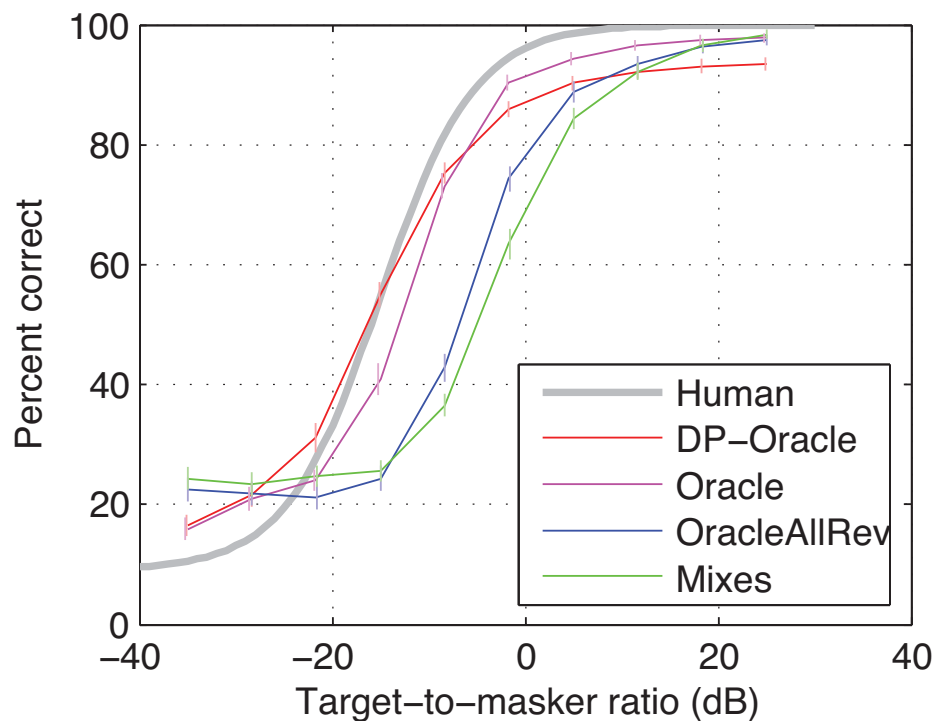


Ground Truth

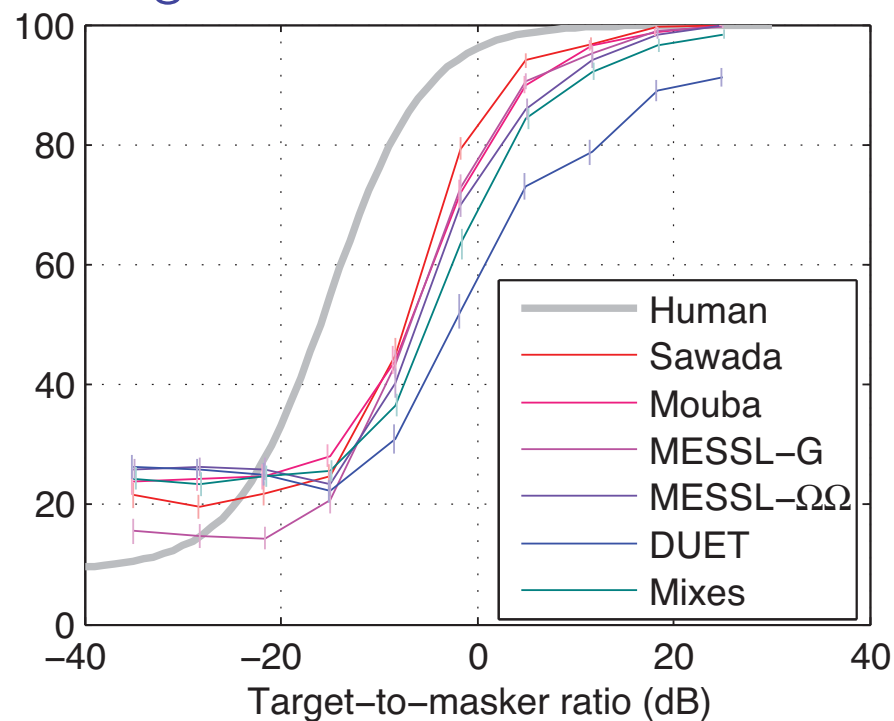
MESSL Results

- Speech recognizer (Digits)

Ground truth masks



Algorithmic masks



3. Combining Spatial + Speech Models

Weiss, Mandel & Ellis '08

- **Interaural** parameters give
 $ILD_i(\omega), ITD_i, \Pr(X(t, \omega) = S_i(t, \omega))$
- **Speech source model** can give
 $\Pr(S_i(t, \omega) \text{ is speech signal})$
- Can combine into one big **EM framework**...

E-step

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



M-step

$$\Theta^{(n+1)} = \operatorname{argmax}_{\Theta} E_{p(u|\Theta^{(n)})} p(x, u|\Theta)$$

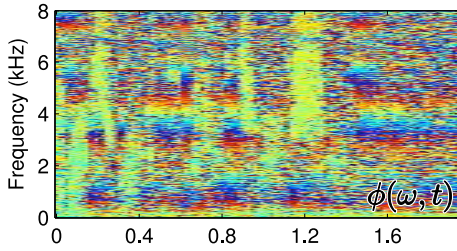
u is: $\Pr(\text{cell from source } i)$
phoneme sequence

Θ is: interaural params
speaker params

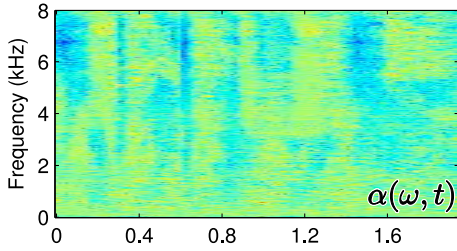
MESSL-SP (Source Prior)

Observations

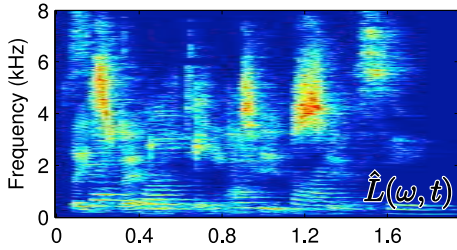
Mixture – IPD



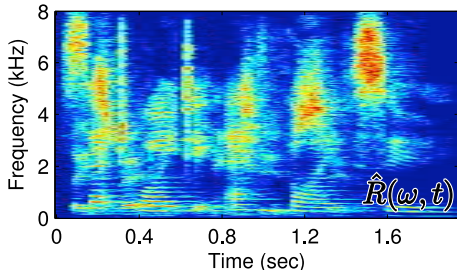
Mixture – ILD



Mixture – left channel

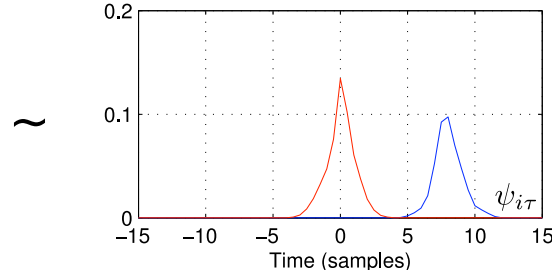


Mixture – right channel

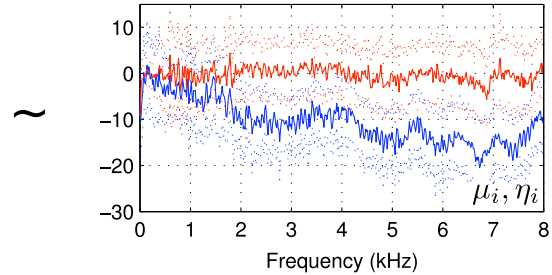


Parameters

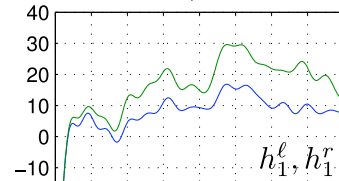
Per-source ITD



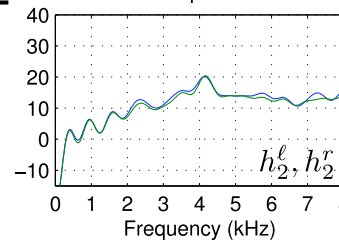
Per-source ILD



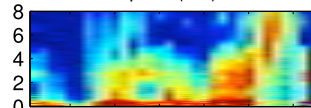
SP channel response – source 1



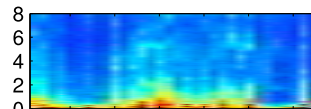
SP channel response – source 2



Source prior (SP) means



SP covars

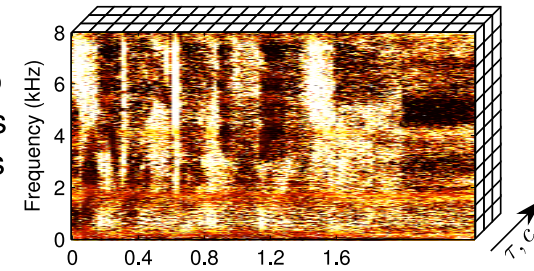


Mixture component

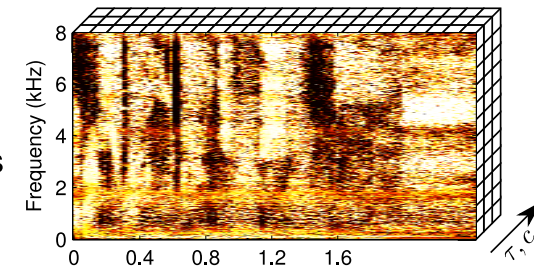
Posteriors

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

Source 1 mask

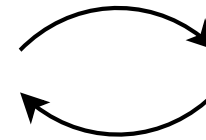


Source 2 mask



E-step

Use parameters to compute posteriors of hidden variables



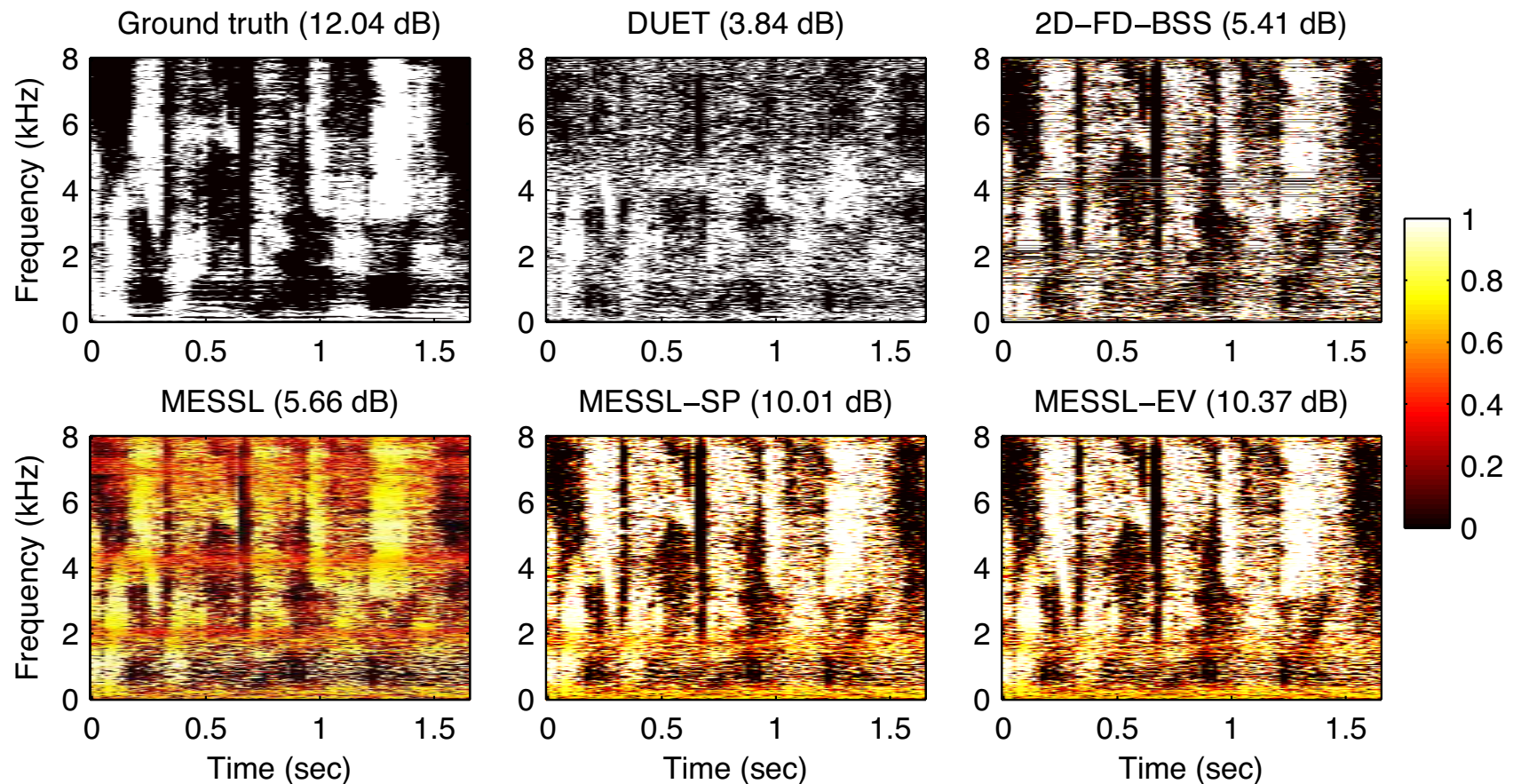
M-step

Use posteriors to update parameters

Separate sources by multiplying mixture by different masks

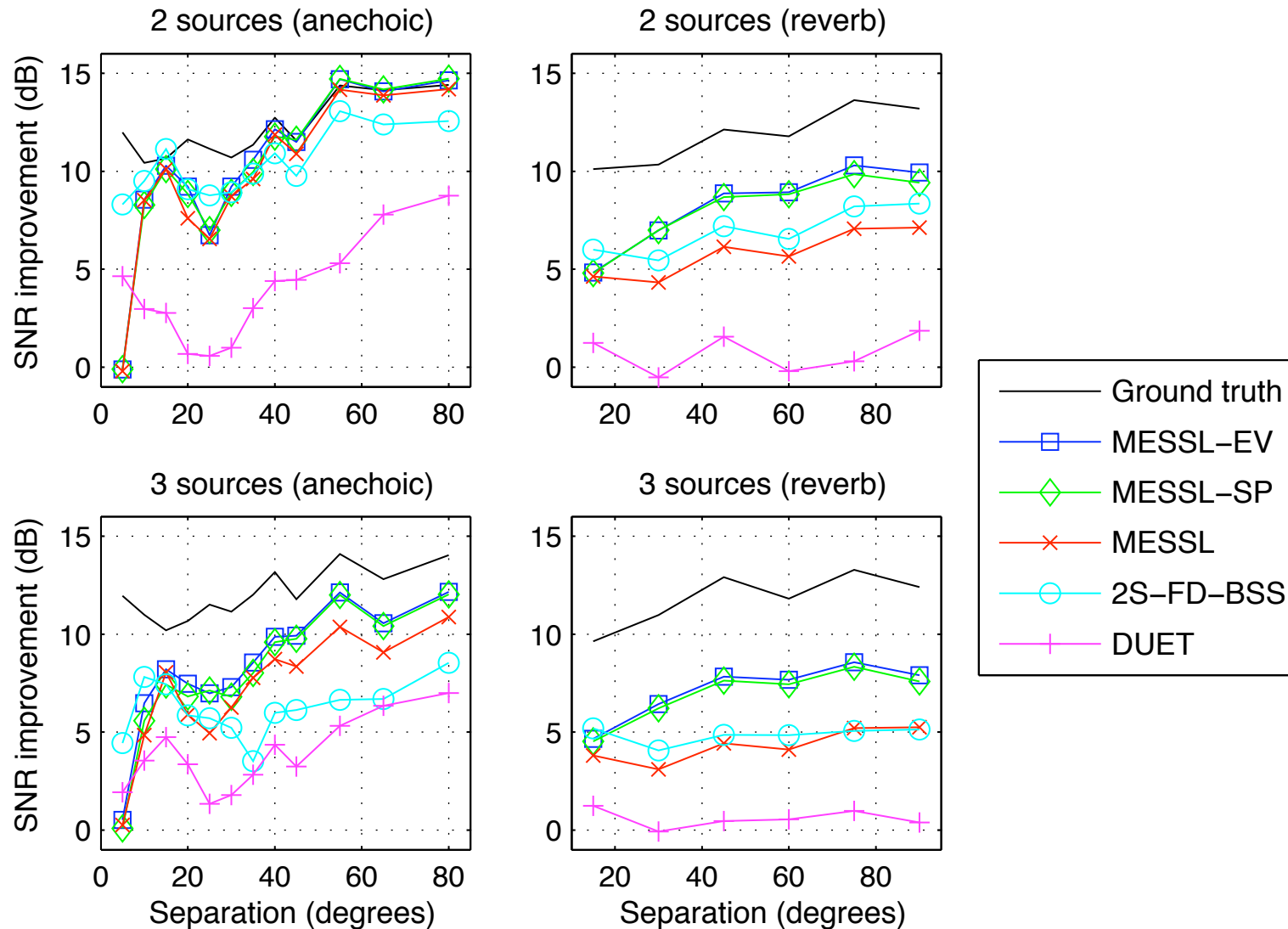
MESSL-SP Results

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



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