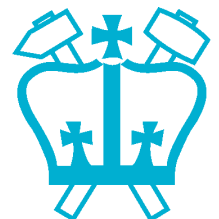# Research in Sound Analysis

## Dan Ellis

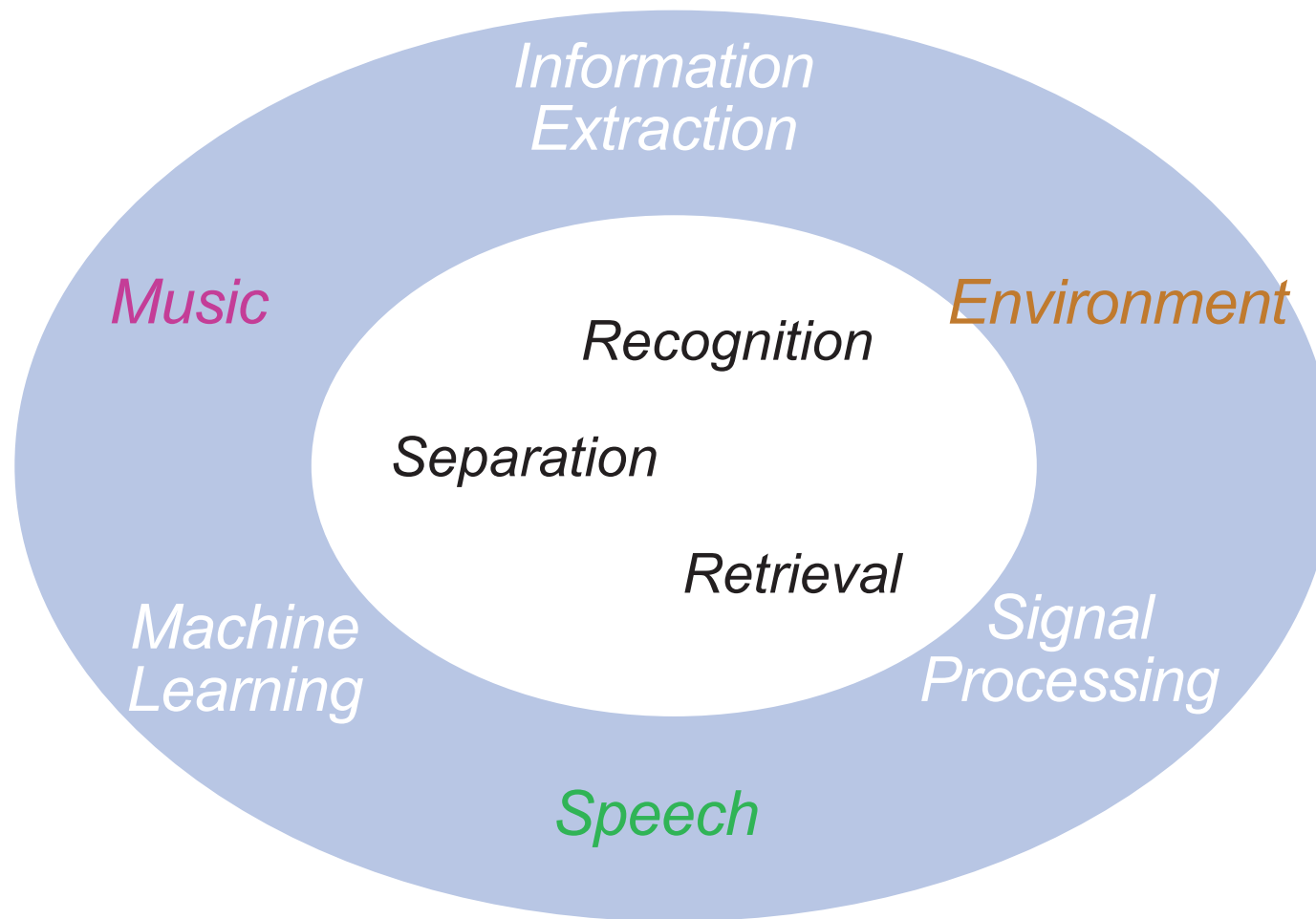Laboratory for **R**ecognition and **O**rganization of **S**peech and **A**udio
Dept. Electrical Eng., Columbia Univ., NY USA

dpwe@ee.columbia.edu

1. Personal and Consumer Audio
2. Musical Cover Song Detection
3. Binaural Source Separation

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

COLUMBIA UNIVERSITY
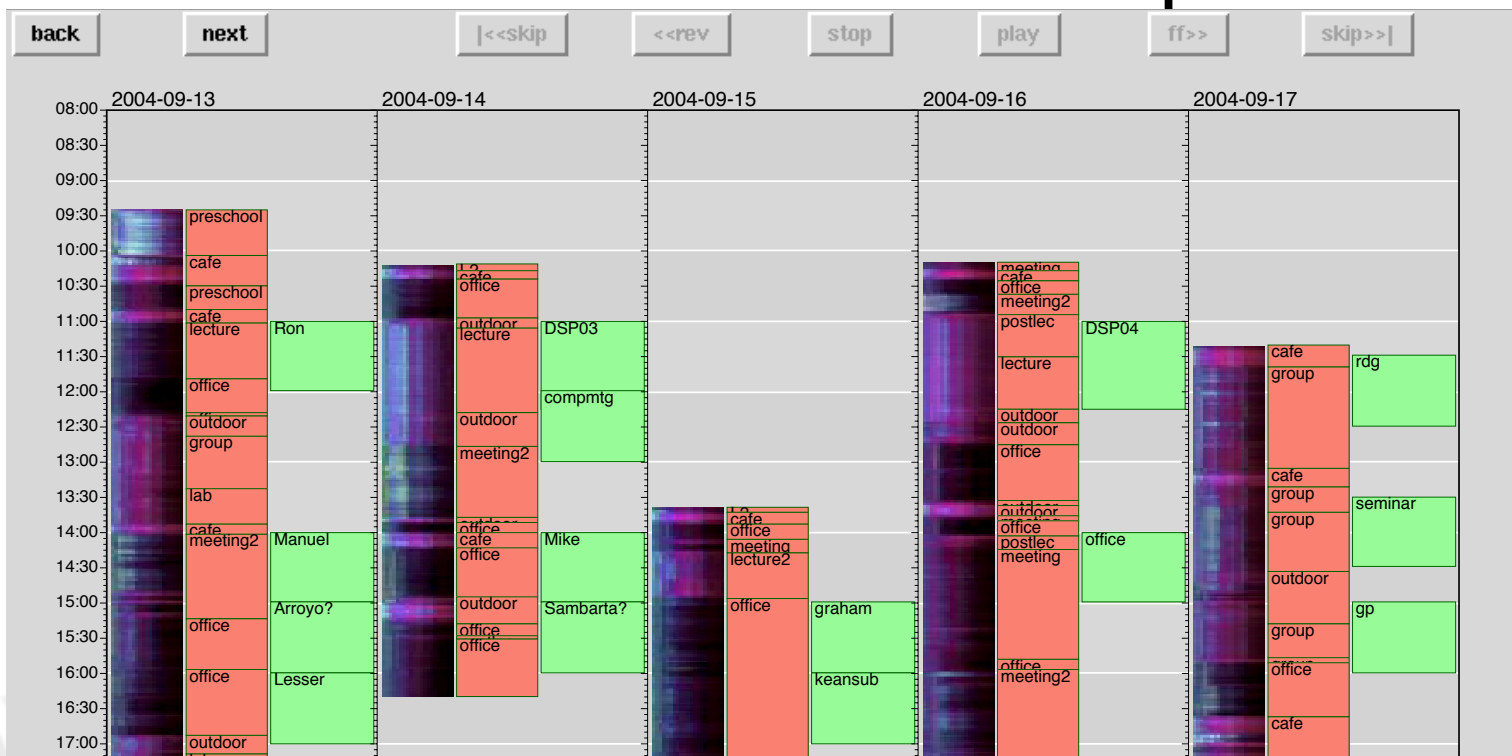IN THE CITY OF NEW YORK

# LabROSA Overview

# I. Personal Audio Archives

- **Easy to record everything you hear**
  - <2GB / week @ 64 kbps

- **Hard to find anything**
  - how to scan?
  - how to visualize?
  - how to index?

- **Need automatic analysis**

- **Need minimal impact**

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Browsing Interface

- Browsing / Diary interface
  - links to other information (diary, email, photos)
  - synchronize with note taking? *(Stifelman & Arons)*
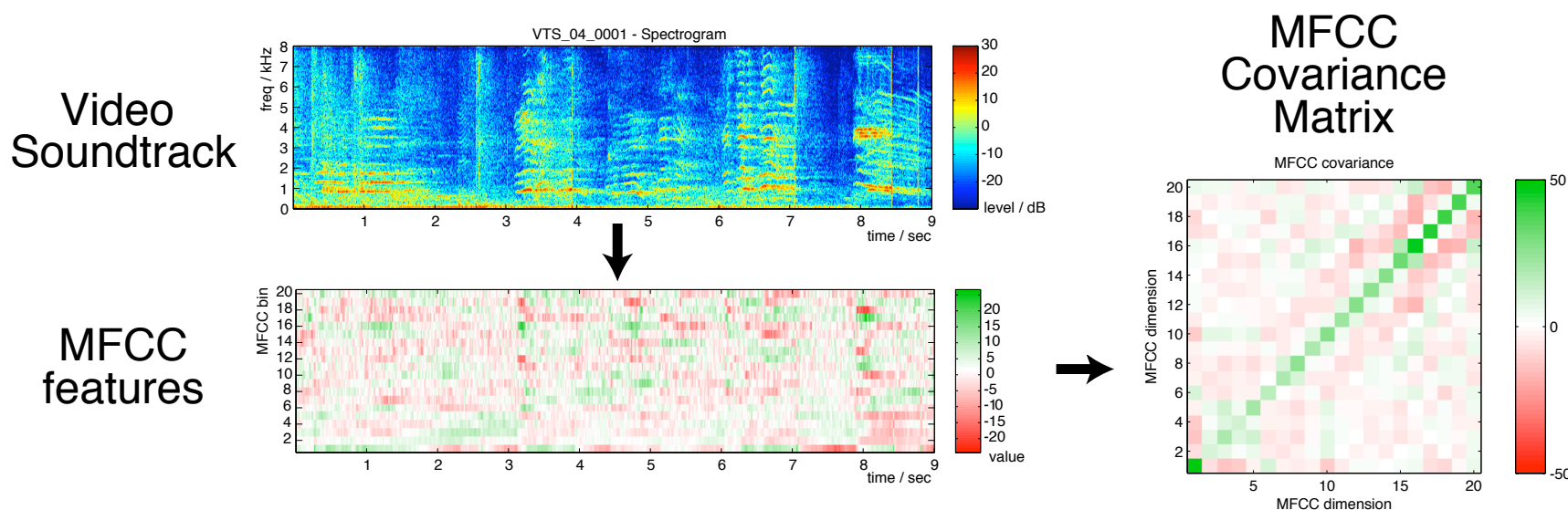  - audio thumbnails

- Release Tools + "how to" for capture

# Consumer Video



- Short video clips as the **evolution of snapshots**
  - 10-60 sec, one location, no editing
  - browsing?

- More information for indexing...
  - video + audio
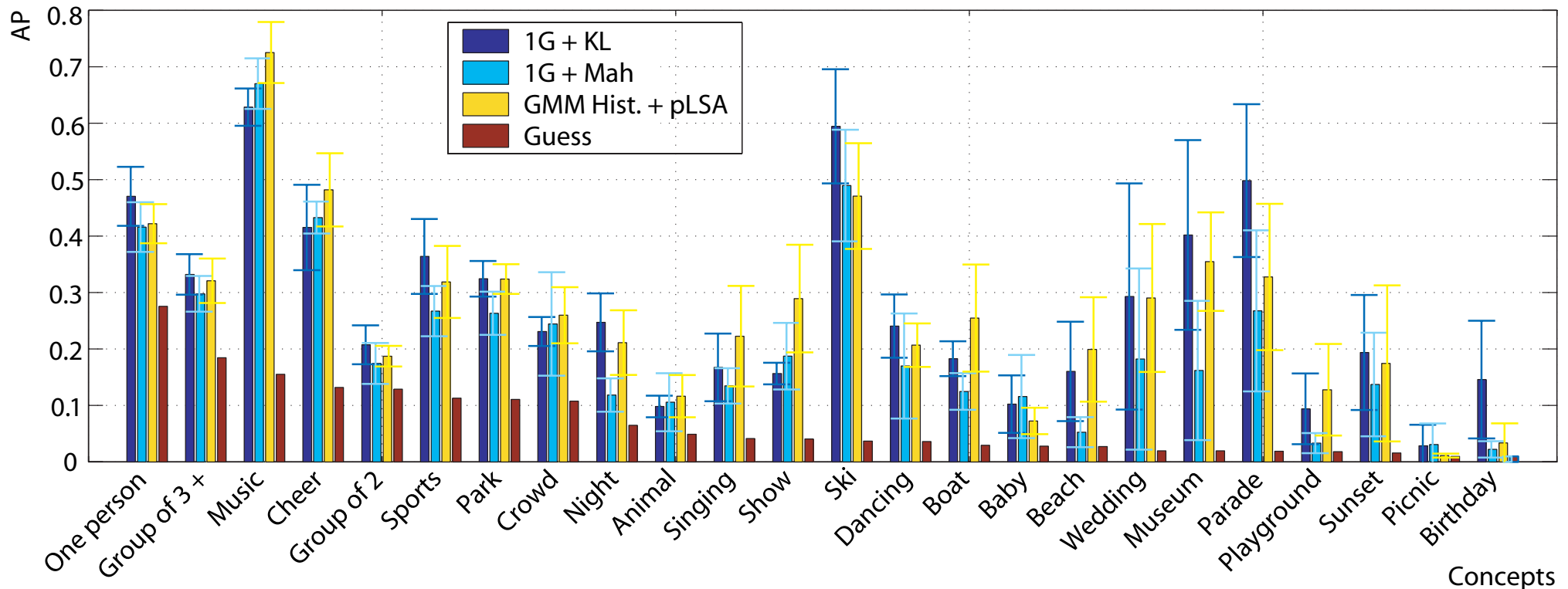  - foreground + background

# MFCC Covariance Representation

- **Each clip/segment → fixed-size statistics**
  - similar to speaker ID and music genre classification
- **Full Covariance matrix of MFCCs**
  - maps the kinds of spectral shapes present



- **Clip-to-clip distances for SVM classifier**
  - by KL or 2nd Gaussian model

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Audio-Only Results

- ## Wide range of results:



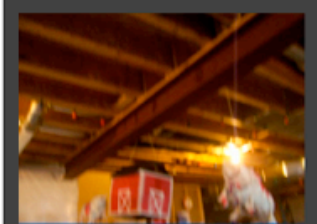- ○ audio (music, ski) vs. non-audio (group, night)
- ○ large AP uncertainty on infrequent classes

# How does it 'feel'?

- Browser impressions: How wrong is wrong?

*Top 8 hits for "Baby"*

# 2. Cover Song Detection: Chroma

- ## Chroma features map spectral energy into one <span style="color:brown">canonical octave</span>
  - i.e. 12 semitone bins

*Piano scale*



Piano chromatic scale



IF chroma

- ## Can resynthesize as <span style="color:green">"Shepard Tones"</span>
  - all octaves at once



12 Shepard tone spectra



Shepard tone resynth

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Beat-Synchronous Chroma Features

- Beat + chroma features / 30ms frames
  → average chroma within each beat
  ○ compact; sufficient?

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Cross-Correlation Matching

- Look inside global cross-correlation to find matching fragments...

  - $\text{xcorr} = \Sigma_t \, \Sigma_f \left( C_1(t,f) \cdot C_2(t,f) \right)$ - view along time



**Let It Be / Beatles (beats 11-441)**

**Let It Be / Nick Cave (beats 13-443)**

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# "The Meaning of Music"

The ultimate goal of this research...

- ## What does music evoke in a listener's mind?
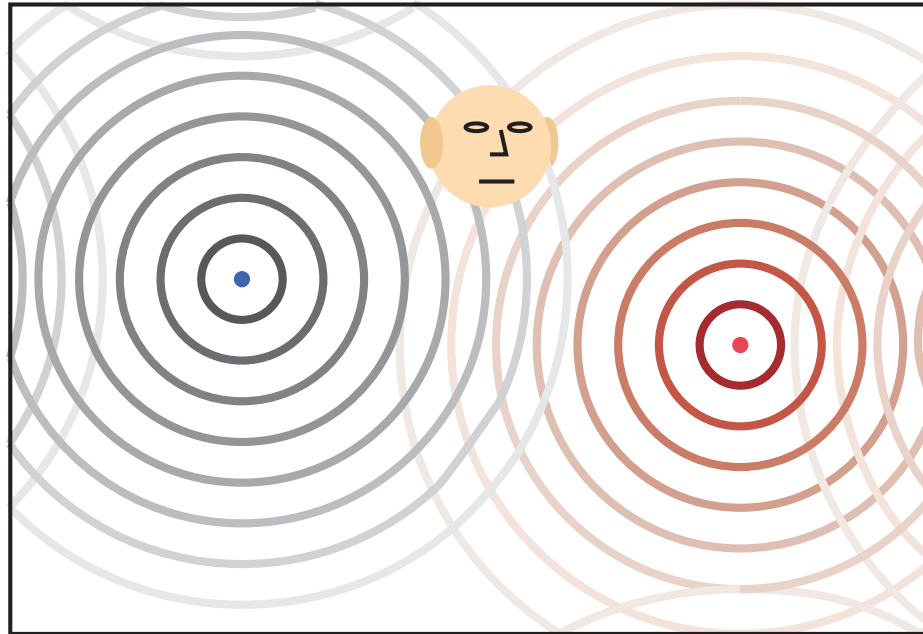  - i.e. "what does it all mean?" (metaphysics?)
  - study with subjective experiments
  - (then build detectors for specific responses ...?)

- ## What phenomena are denoted by "music"?
  - i.e. delineate the "set of all music"
  - (the ultimate music/nonmusic classifier?)

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# 3. Binaural Source Separation

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



- Tasks:
  - identify positions of sources (and number?)
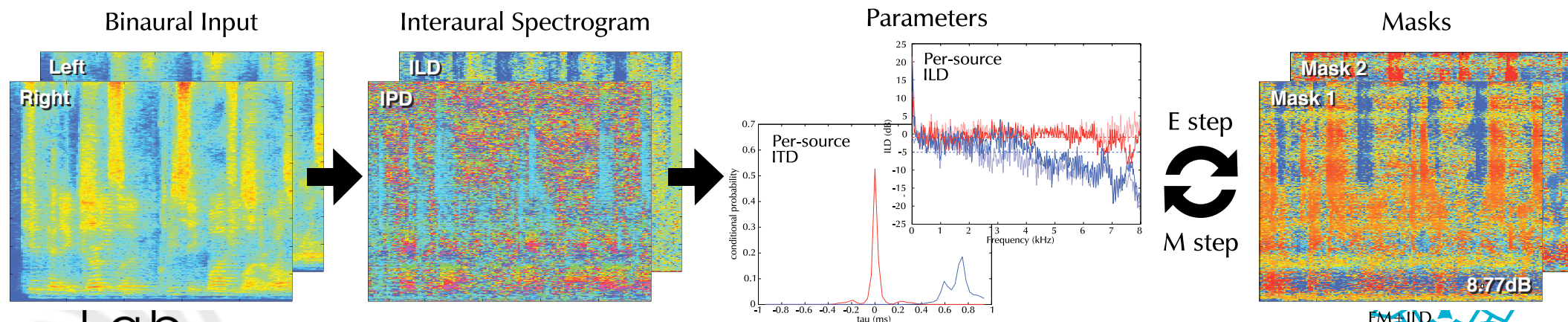  - recover source signals

# Spatial Estimation in Reverb

*Mandel & Ellis '07*

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

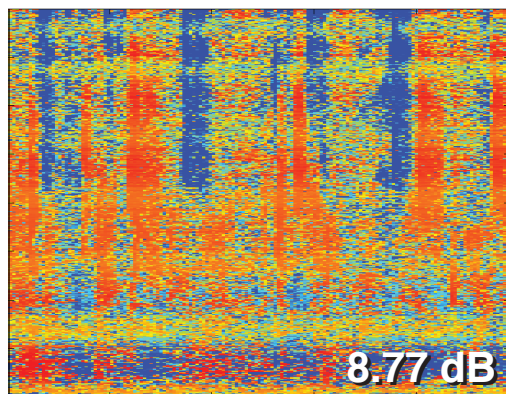$$\frac{L(\omega, t)}{R(\omega, t)} = a(\omega)e^{j\omega\tau}N(\omega, t)$$

  - converge via EM to $a()$, $\tau$ for each source
  - mask is $\Pr(X(t,\omega)$ dominated by source $i)$



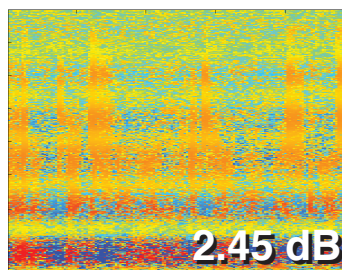Binaural Input · Interaural Spectrogram · Parameters · Masks

# Spatial Estimation Results

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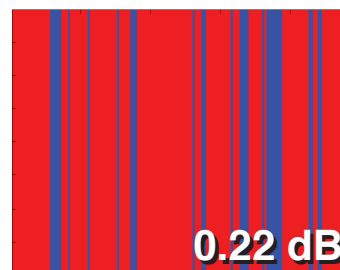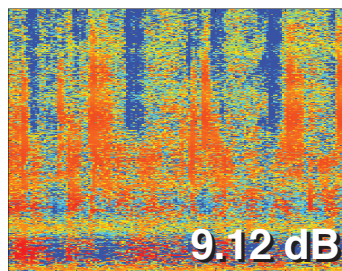  - tradeoff between constraints & noisiness
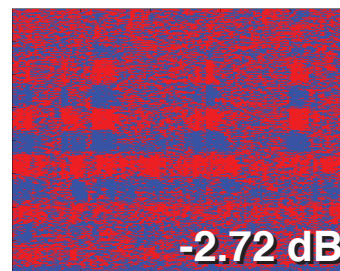


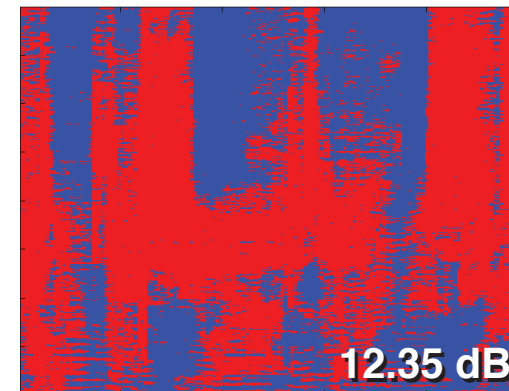EM+ILD — 8.77 dB

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

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Conclusions

- **LabROSA**
  - information from sound ...
  - ... via signal processing and machine learning
- **Environmental Sound**
- **Music Audio**
- **Source Separation**
- **Speech, models, dolphins...**

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK