# Using Sound Source Models to Organize Mixtures

## Dan Ellis

Laboratory for Recognition and Organization of Speech and Audio
Dept. Electrical Eng., Columbia Univ., NY USA
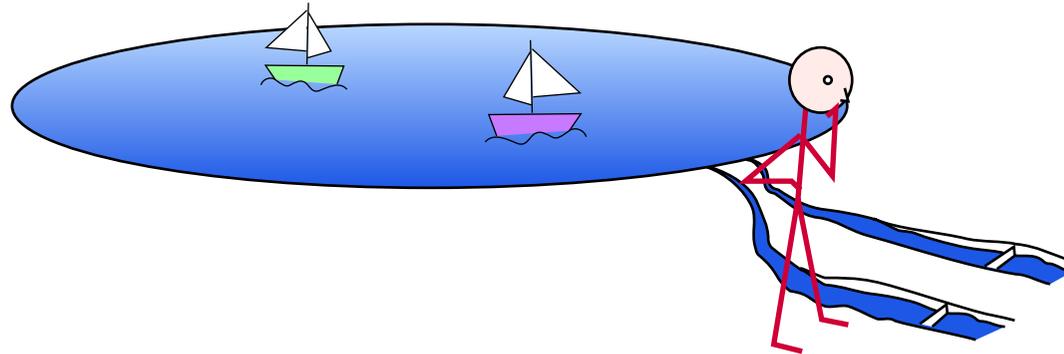
dpwe@ee.columbia.edu                    http://labrosa.ee.columbia.edu/

1.  Mixtures and Models
2.  Human Sound Organization
3.  Machine Sound Organization
4.  Ambient Sounds

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

COLUMBIA UNIVERSITY
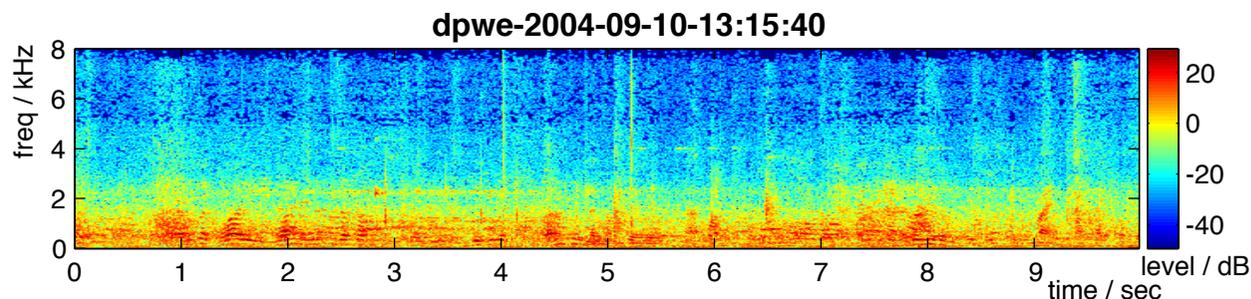IN THE CITY OF NEW YORK

# The Problem of Mixtures



*"Imagine two narrow channels dug up from the edge of a lake, with handkerchiefs stretched across each one. Looking only at the motion of the handkerchiefs, you are to answer questions such as: How many boats are there on the lake and where are they?"*  (after Bregman'90)

- ## Received waveform is a mixture
  - ○ 2 sensors, N sources - underconstrained

- ## Undoing mixtures: hearing's primary goal?
  - ○ .. by any means available

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Sound Organization Scenarios

- **Interactive voice systems**
  - human-level understanding is expected
- **Speech prostheses**
  - crowds: #1 complaint of hearing aid users
- **Archive analysis**
  - identifying and isolating sound events



dpwe-2004-09-10-13:15:40

pa-2004-09-10-131540.wav

- Unmixing/remixing/enhancement...

LabROSA
Laboratory for the Recognition and
Organization of Speech and Audio

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# How Can We Separate?

- **By between-sensor differences** (spatial cues)
  - 'steer a null' onto a compact interfering source
  - the filtering/signal processing paradigm

- **By finding a 'separable representation'**
  - spectral?  sources are broadband but sparse
  - periodicity?  maybe – for pitched sounds
  - something more signal-specific...

- **By inference (based on knowledge/models)**
  - acoustic sources are redundant
    - → use part to guess the remainder
    - - limited possible solutions

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

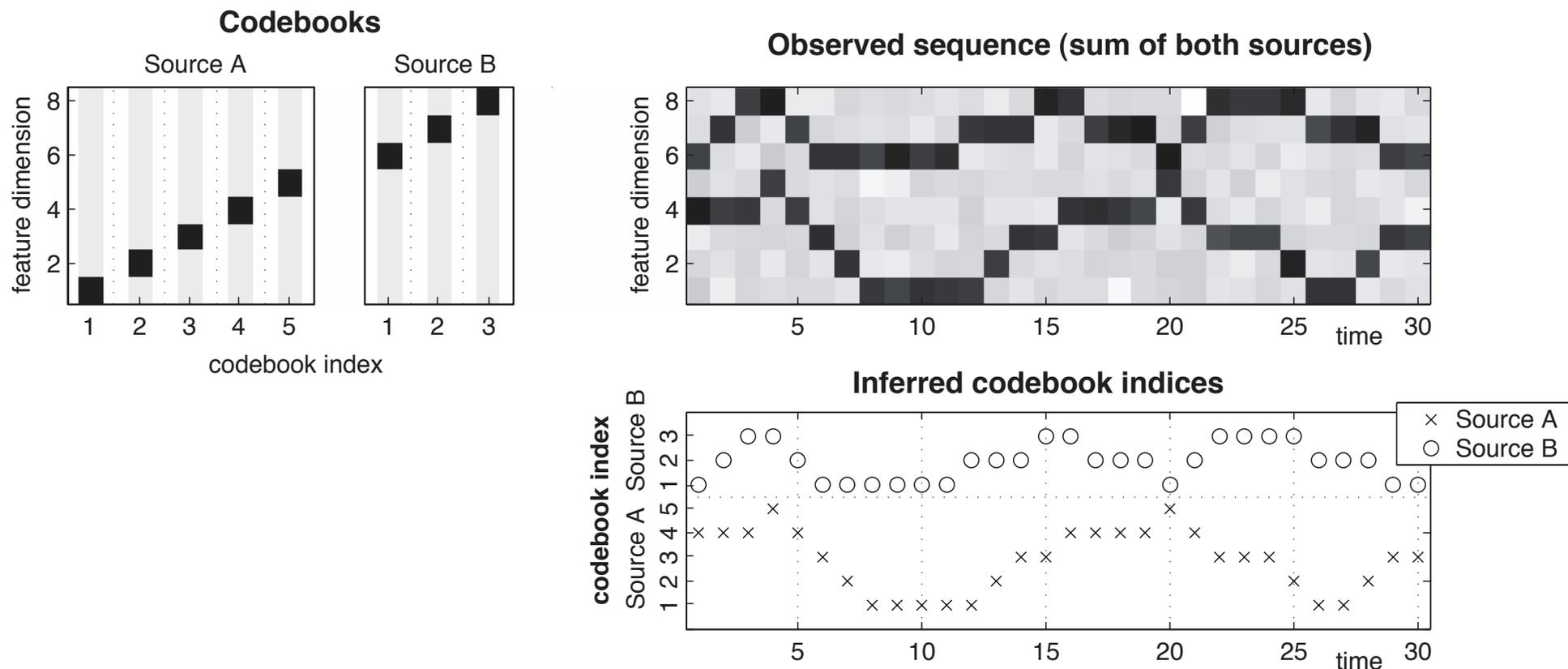COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Separation vs. Inference

- **Ideal** separation is rarely possible
  - i.e. no projection can completely remove overlaps
- Overlaps → Ambiguity
  - scene analysis = find "most reasonable" explanation
- Ambiguity can be expressed probabilistically
  - i.e. posteriors of sources $\{S_i\}$ given observations $X$:

$$P(\{S_i\}|\ X) \propto P(X\ |\{S_i\})\ P(\{S_i\})$$

$$\text{combination physics} \quad \text{source models}$$

- Better **source models** → better **inference**

  - .. learn from examples?

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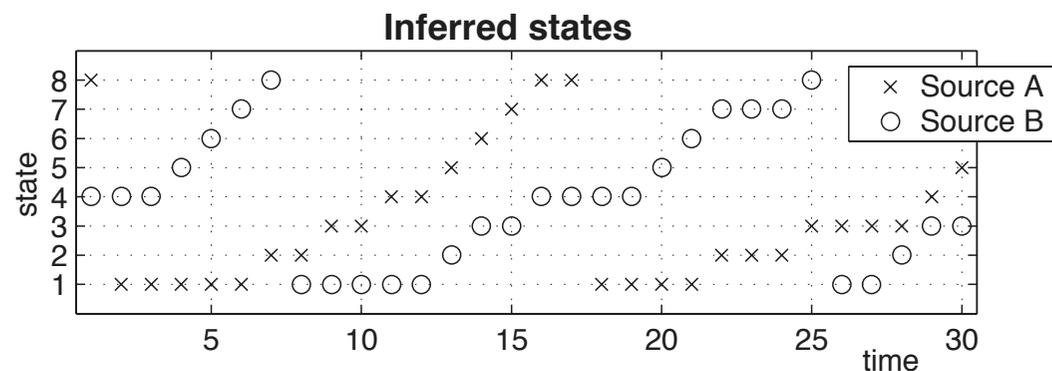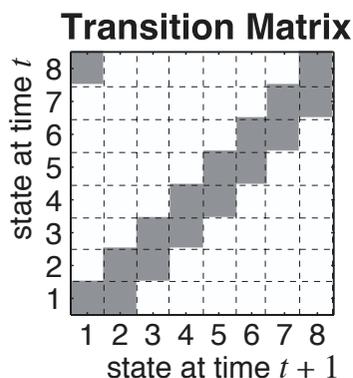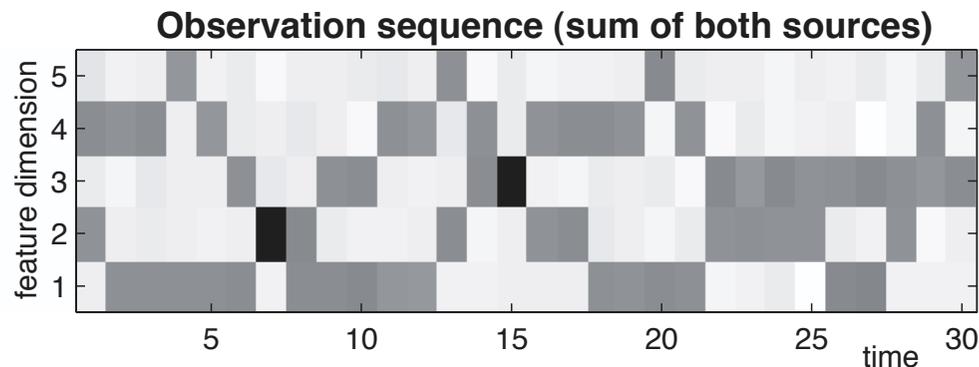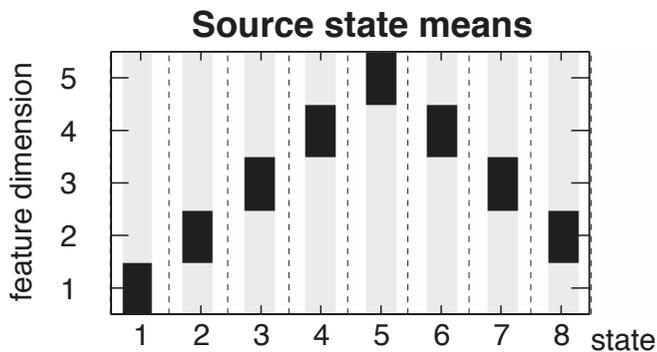Organization of Speech and Audio

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# A Simple Example

- Source models are codebooks from separate subspaces

# A Slightly Less Simple Example

- Sources with Markov transitions

**Source state means**

**Observation sequence (sum of both sources)**

**State diagram**

**Transition Matrix**

**Inferred states**

× Source A
○ Source B

Lab ROSA

Laboratory for the Recognition and
Organization of Speech and Audio

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

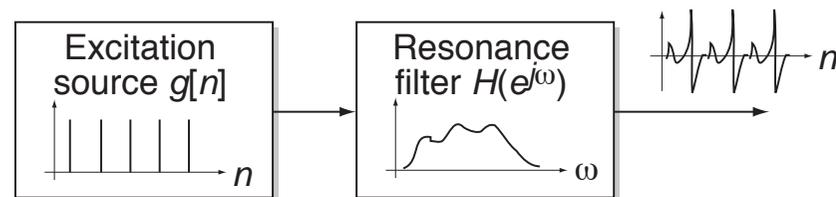# What is a Source Model?

- **Source Model** describes signal behavior
  - encapsulates constraints on form of signal
  - (any such constraint can be seen as a model...)

- A model has **parameters**
  - model + parameters
    - → instance



Excitation source $g[n]$ → Resonance filter $H(e^{j\omega})$ → $n$

- What is *not* a source model?
  - detail not provided in instance
    e.g. using phase from original mixture
  - constraints on interaction between sources
    e.g. independence, clustering attributes

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Outline

1. Mixtures and Models
2. Human Sound Organization
   - Auditory Scene Analysis
   - Using source characteristics
   - Illusions
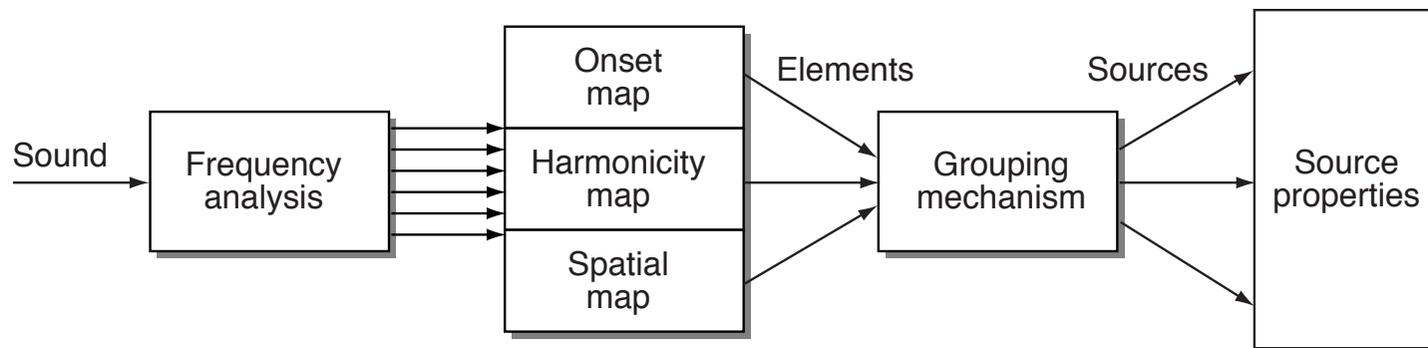3. Machine Sound Organization
4. Ambient Sounds

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Auditory Scene Analysis

*Bregman'90*
*Darwin & Carlyon'95*

- **How do people analyze sound mixtures?**
  - ○ break mixture into small elements (in time-freq)
  - ○ elements are grouped in to sources using cues
  - ○ sources have aggregate attributes
- **Grouping rules (Darwin, Carlyon, ...):**
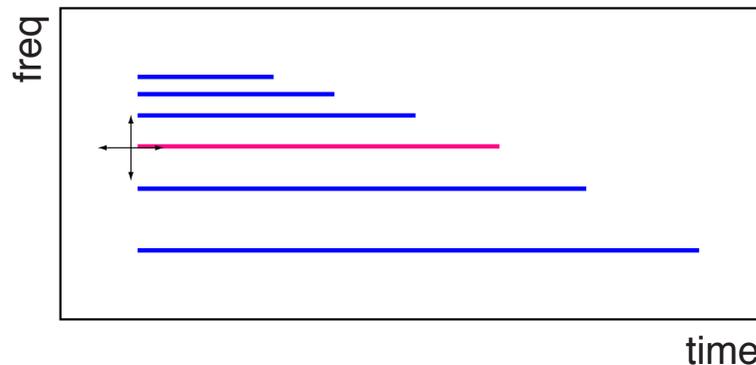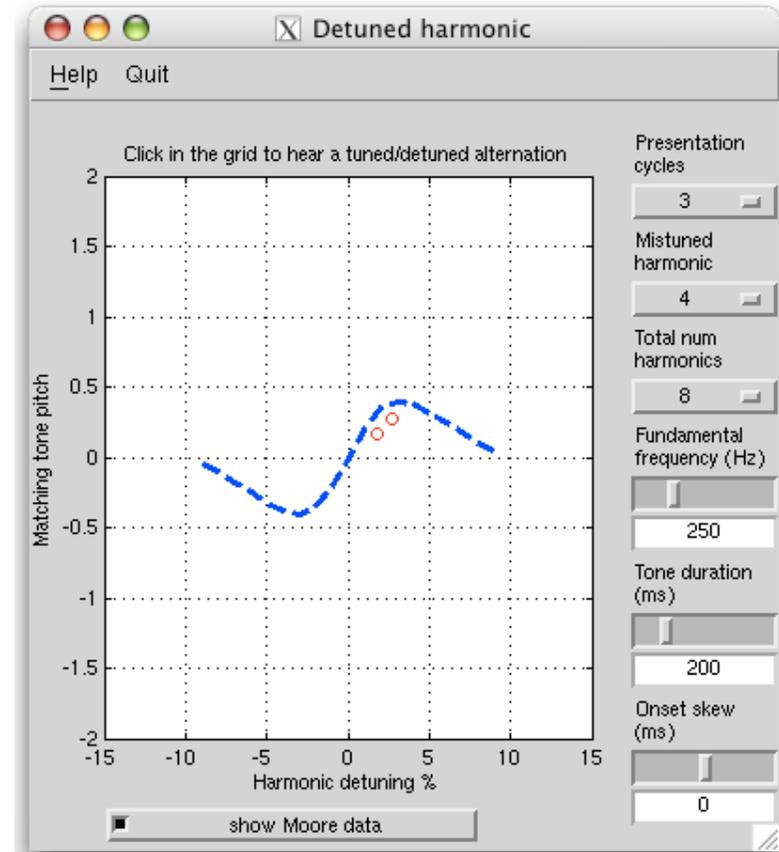  - ○ cues: common onset/modulation, harmonicity, ...



*(after Darwin 1996)*

- **Also learned "schema" (for speech etc.)**

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Perceiving Sources

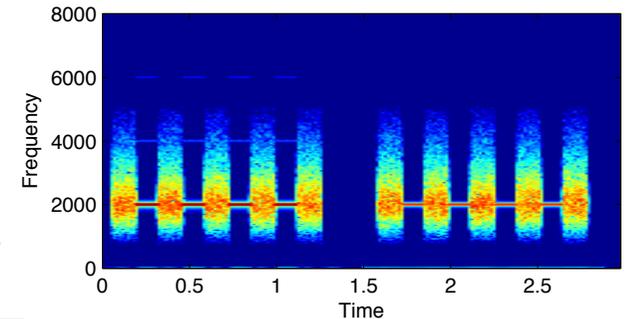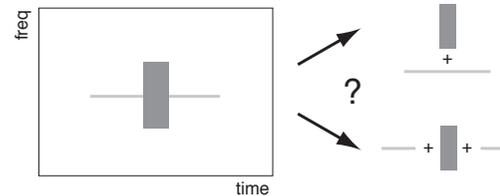- **Harmonics** distinct in ear, but perceived as one source ("fused"):



- depends on common onset
- depends on harmonics

- **Experimental techniques**
  - ask subjects "how many"
  - match attributes e.g. pitch, vowel identity
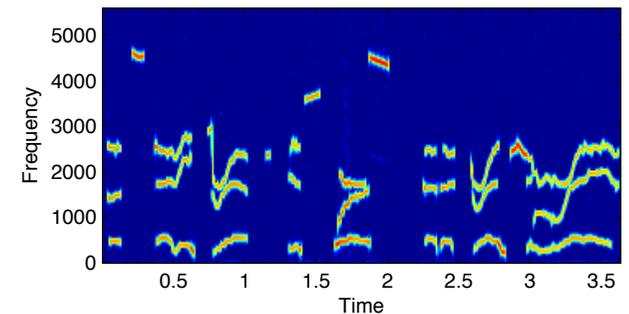  - brain recordings (EEG "mismatch negativity")

# Auditory "Illusions"

- ## How do we explain illusions?
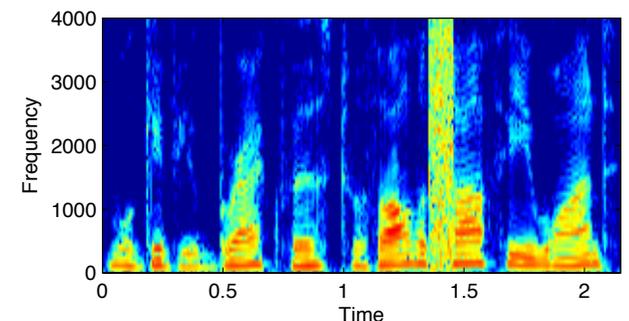  - ○ pulsation threshold

  - ○ sinewave speech

  - ○ phonemic restoration

- ## Something is providing the missing (illusory) pieces ... source models

LabROSA
Laboratory for the Recognition and
Organization of Speech and Audio

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Human Speech Separation

*Brungart et al.'02*

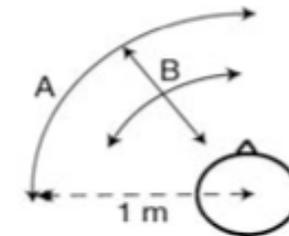- **Task:** Coordinate Response Measure
  - "Ready Baron go to green eight now"
  - 256 variants, 16 speakers
  - correct = color and number for "Baron"

- **Accuracy as a function of spatial separation:**



  - A, B same speaker
  - Range effect

# Separation by Vocal Differences

- ## CRM varying the level and voice character



*(same spatial location)*

- energetic vs. informational masking
- more than pitch .. source models

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Outline

1. Mixtures and Models
2. Human Sound Organization
3. **Machine Sound Organization**
   - Computational Auditory Scene Analysis
   - Dictionary Source Models
4. Ambient Sounds

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Source Model Issues

- Domain
  - parsimonious expression of constraints
  - nice combination physics
- Tractability
  - size of search space
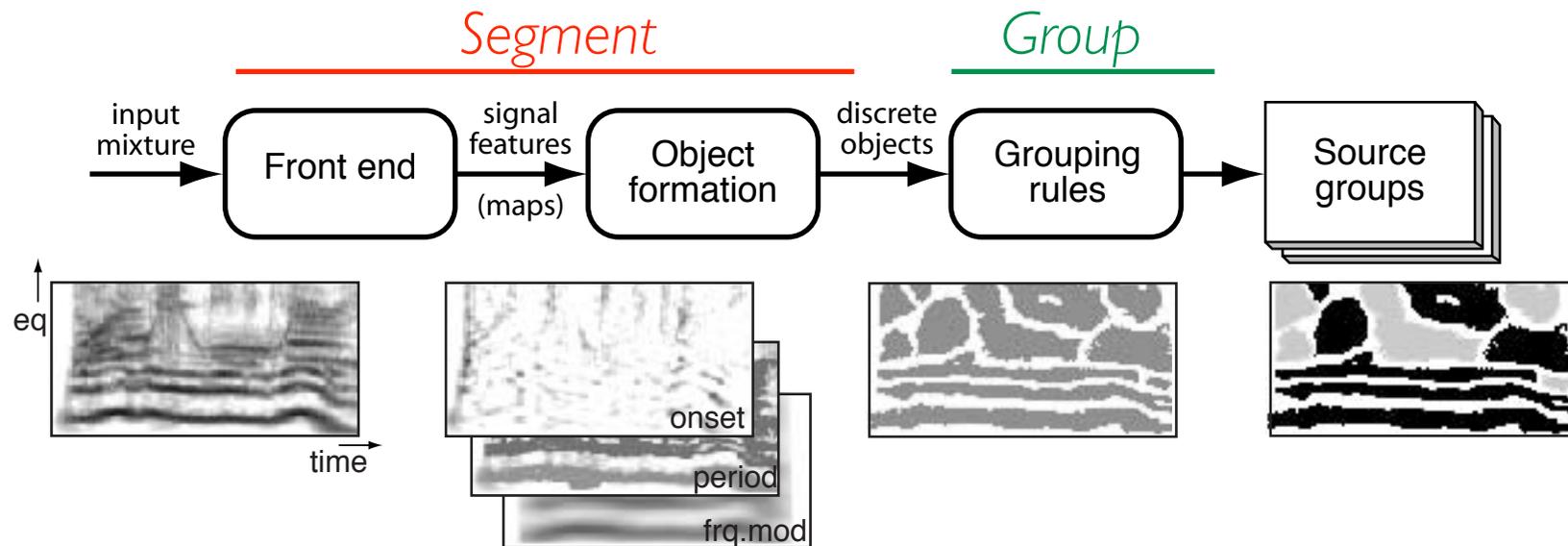  - tricks to speed search/inference
- Acquisition
  - hand-designed vs. learned
  - static vs. short-term
- Factorization
  - independent aspects
  - hierarchy & specificity

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Computational Auditory Scene Analysis

*Brown & Cooke'94*
*Okuno et al.'99*
*Hu & Wang'04 ...*

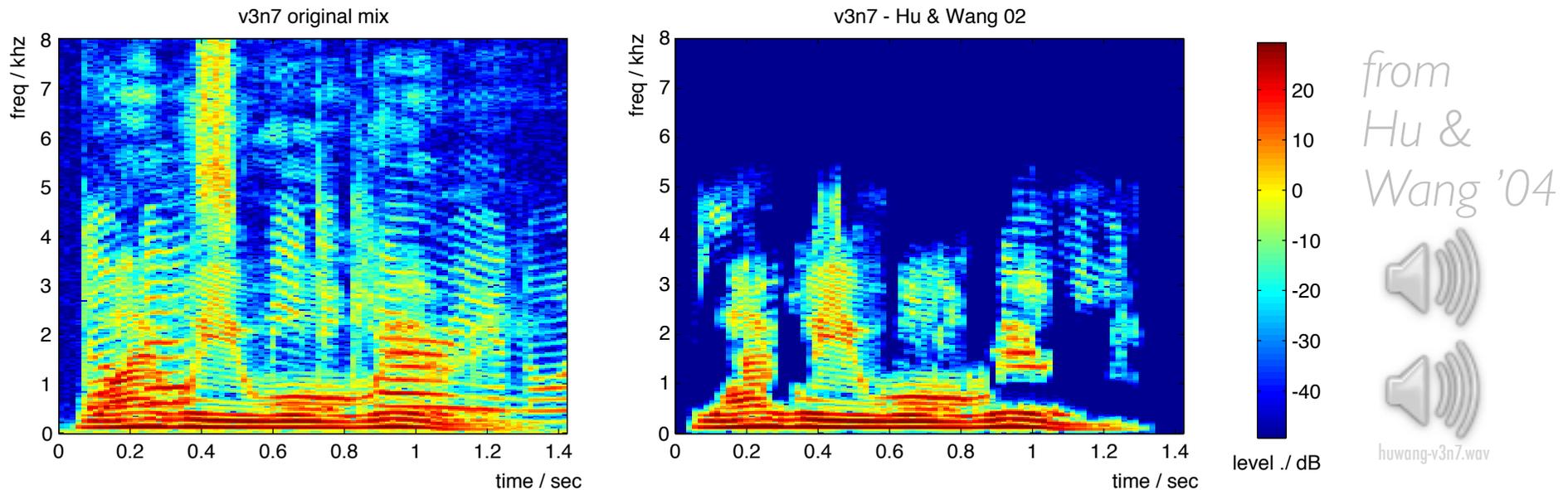- ## Central idea:
  Segment time-frequency into sources based on perceptual grouping cues



- ... principal cue is harmonicity

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# CASA limitations

- ## Limitations of T-F masking
  - ○ cannot undo overlaps – leaves gaps



v3n7 original mix

v3n7 - Hu & Wang 02

from Hu & Wang '04

huwang-v3n7.wav

- ## Typically driven by local features
  - ○ limited model scope → no inference or illusions
- ## Processing hand-defined, not learned

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Can Models Do CASA?

- **Source models** can learn harmonicity, onset
  - ○ ... to subsume rules/representations of CASA



VQ800 Codebook - Linear distortion measure

  - ○ can capture spatial info too *[Pearlmutter & Zador'04]*

- **Can also capture sequential structure**
  - ○ e.g. consonants follow vowels
  - ○ ... like people do?

- **But: need source-specific models**
  **... for every possible source**
  - ○ use model adaptation? *[Ozerov et al. 2005]*

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Separation or Description?

- **Are isolated waveforms required?**
  - clearly sufficient, but may not be necessary
  - not part of perceptual source separation!
- **Integrate separation with application?**
  - e.g. speech recognition



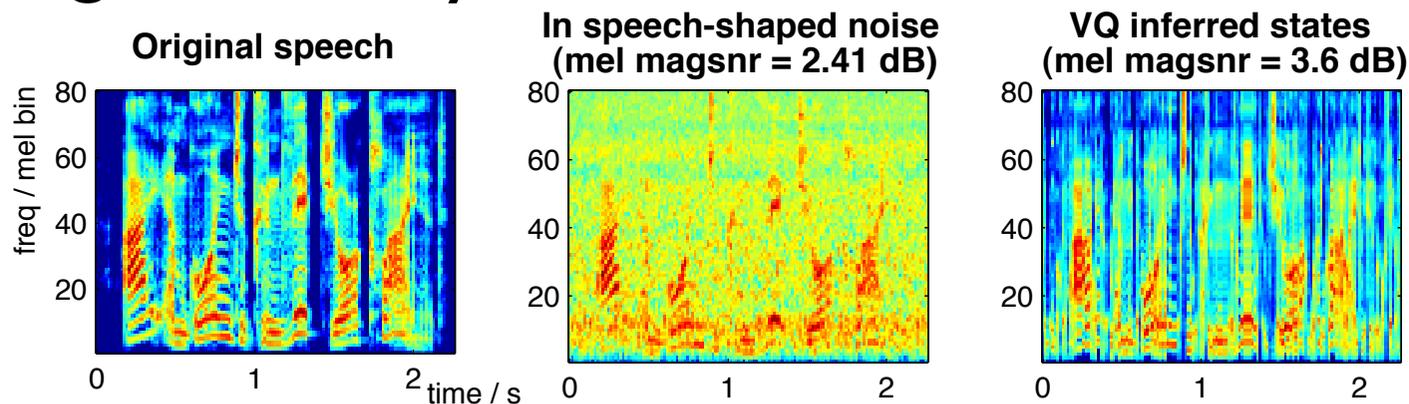  - words output = abstract description of signal

# Dictionary 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...
- ○ different domains for combining $\mathbf{c}$ and defining $\Sigma$

- E.g. stationary noise:



**Original speech**

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

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

freq / mel bin

time / s

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Speech Recognition Models

- **Cooke & Lee Speech Separation Challenge**
  - 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"

- **Decode with Factorial HMM**
  - i.e. two state sequences, one model for each voice
  - exploit sequence constraints
  - exploit speaker differences

- **IBM "superhuman" system** *Kristjansson, Hershey et al. '06*
  - fewer errors than people for same speaker, level
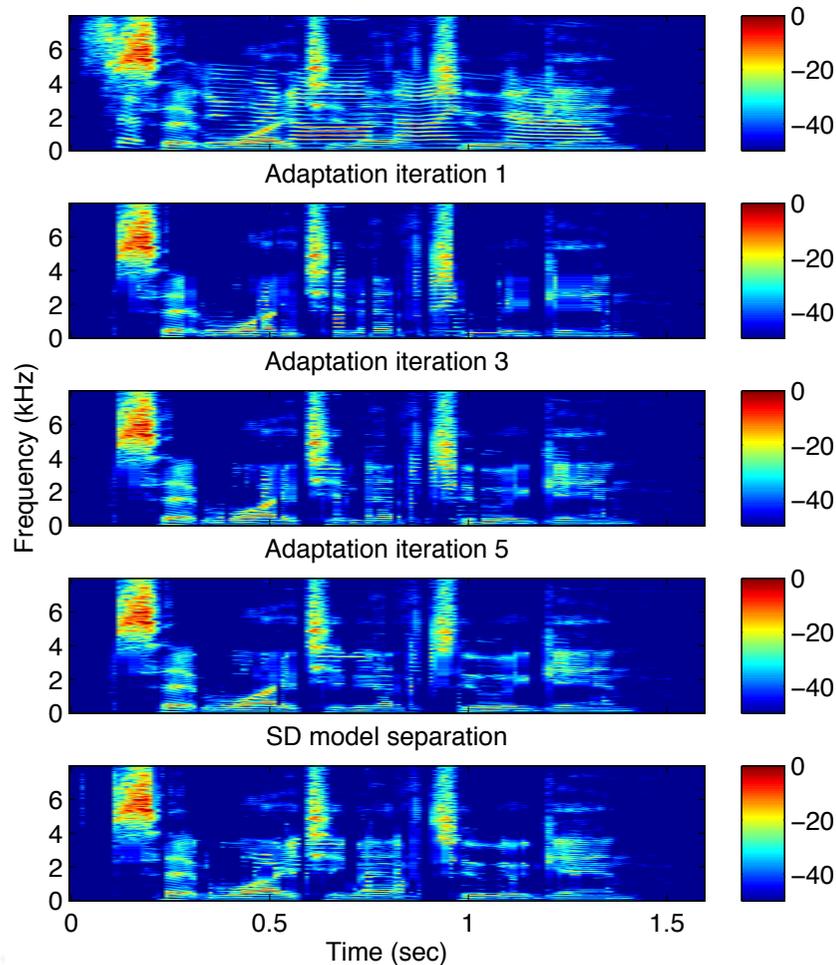  - exploits known speakers, limited grammar

# Speaker-Adapted (SA) Models

*Ron Weiss*

- Factorial HMM needs distinct speakers

Mixture: t32_swil2a_m18_sbar9n

Adaptation iteration 1

Adaptation iteration 3

Adaptation iteration 5

SD model separation

Frequency (kHz)

Time (sec)

- use "eigenvoice" speaker space
- iterate estimating voice & separating speech
- performs midway between speaker-independent (SI) and speaker-dependent (SD)

*SI*

*SA*

*SD*

Diff Gender

acc %

Oracle  SD  SA  SI  Baseline

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# (Pitch) Factored Dictionaries

- ## Separate representations for "source" (pitch) and "filter"
  - $NM$ codewords from $N+M$ entries
  - but: overgeneration...

- ## Faster search
  - direct extraction of pitches
  - immediate separation of (most of) spectra

# Outline

1. Mixtures & Models
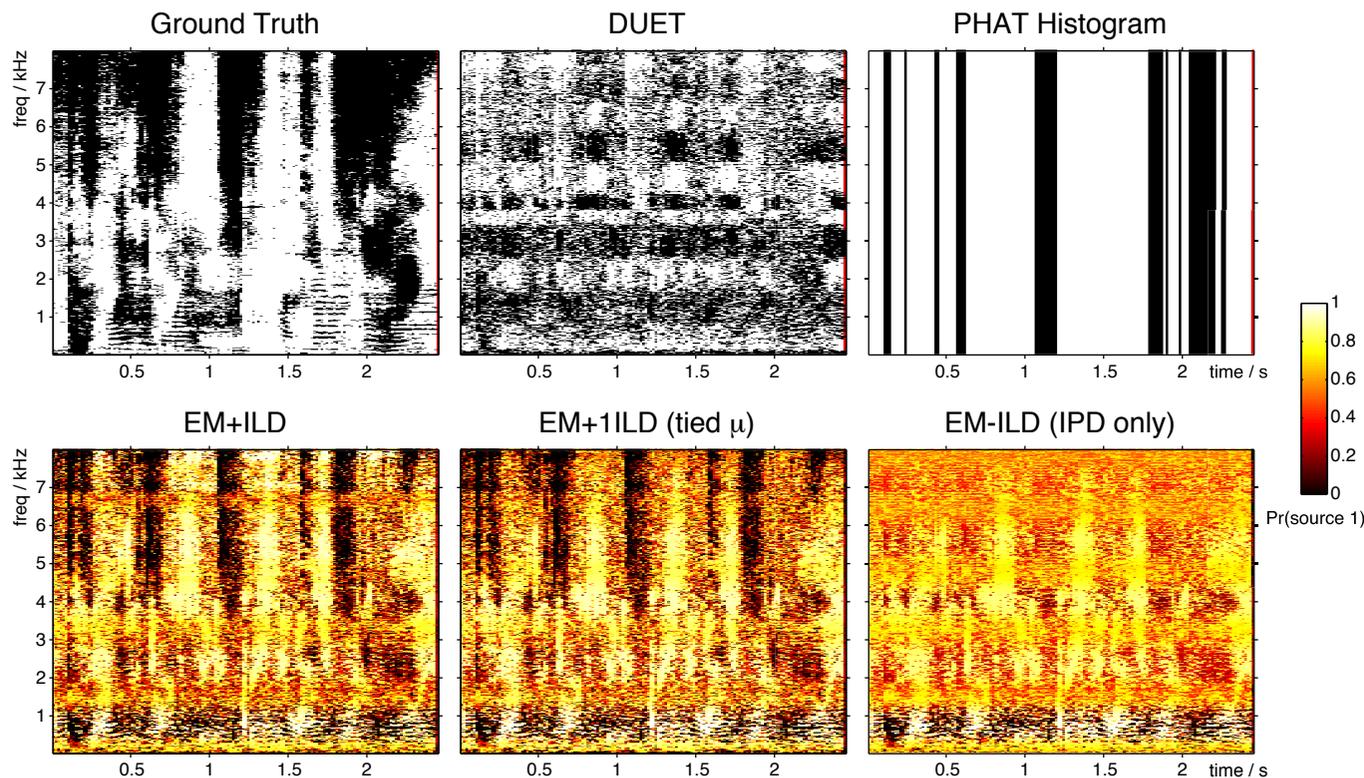2. Human Sound Organization
3. Machine Sound Organization
4. **Ambient Sounds**
   - binaural separation
   - "personal audio" analysis

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Binaural Localization by EM

*Mike Mandel, NIPS'06*

- 2 or 3 sources in reverberation

- Iteratively estimate ILD, IPD
  - initialize from PHAT ITD histogram
  - output is soft TF mask



Ground Truth · DUET · PHAT Histogram · EM+ILD · EM+1ILD (tied μ) · EM-ILD (IPD only)

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# "Personal Audio" Archives

- Continuous recordings with MP3 player
- Segment / cluster "episodes"
  - .. by statistics of ~10 s segments
  - .. for browsing interface

# Personal Audio Speech Detection

*Keansub Lee, Interspeech'06*

- **Pitch is last speech cue to disappear**
  - noise robust pitch tracker for voice detection
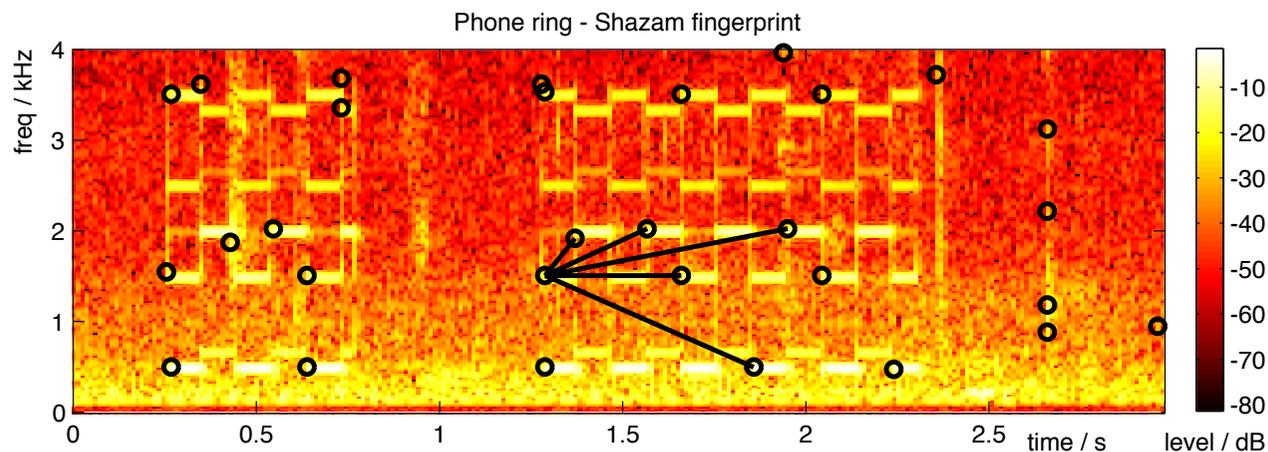  - biggest problem was periodic noise (air conditioning)

Personal Audio - Speech + Noise

Pitch Track + Speaker Active Ground Truth

Lab ROSA

Laboratory for the Recognition and
Organization of Speech and Audio

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Repeating Events in Personal Audio

*Jim Ogle, ICASSP'07*

- "Unsupervised" feature to help browsing
- Full NxN search is very expensive
  - use Shazam fingerprint hashes to find repeats



Phone ring - Shazam fingerprint

  - only works for exact repeats (alarms, jingles)
- O(N) scan for repeats
  - fixed-size hash table
  - multiple common hashes → confident match

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Summary & Conclusions

- **Listeners do well separating sound mixtures**
  - using signal cues (location, periodicity)
  - using source-property variations
- **Machines do less well**
  - difficult to apply enough constraints
  - need to exploit signal detail
- **Models capture constraints**
  - learn from the real world
  - adapt to sources
- **Separation feasible only sometimes**
  - describing source properties is easier

LabROSA
Laboratory for the Recognition and
Organization of Speech and Audio

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK