# Sound, Mixtures, and Learning: LabROSA overview

Dan Ellis
<dpwe@ee.columbia.edu>

Laboratory for **R**ecognition and **O**rganization of **S**peech and **A**udio
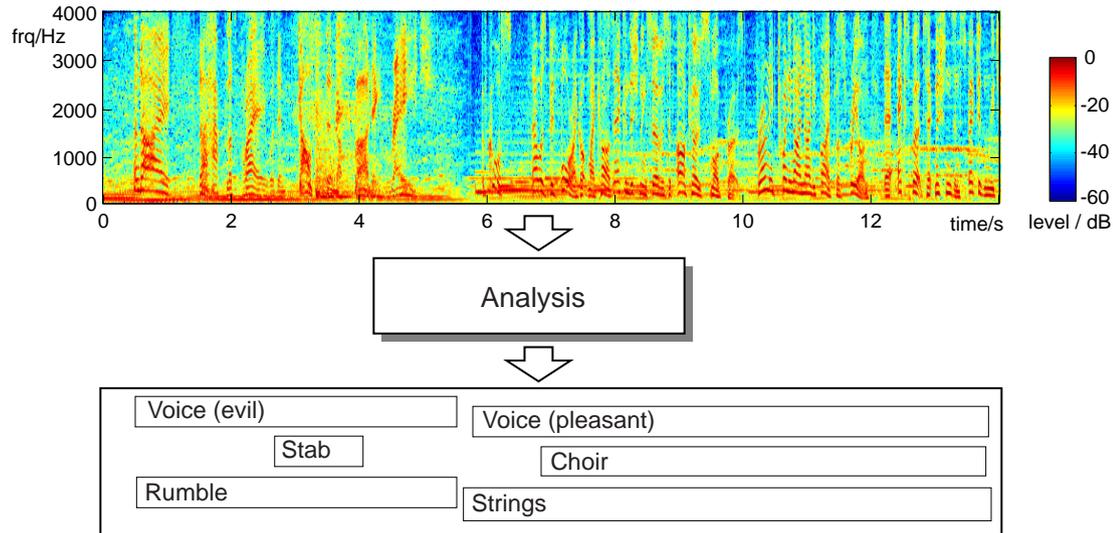Electrical Engineering Dept., Columbia University, New York
http://labrosa.ee.columbia.edu/

## Outline

**1** **Auditory Scene Analysis**

**2** **Speech Recognition & Mixtures**

**3** **Music Analysis & Similarity**

**4** **General Sound Organization**

**5** **Future Work**

Lab
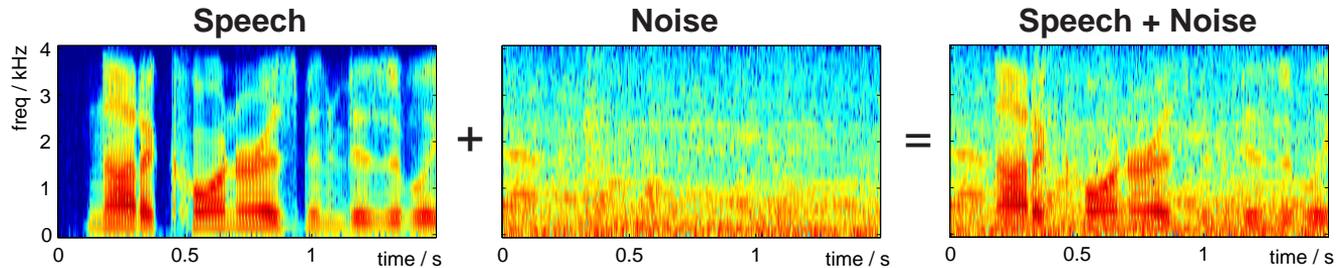ROSA
Laboratory for the Recognition and

# Auditory Scene Analysis



- ***Auditory Scene Analysis*: describing a complex sound in terms of high-level sources/events**
  - ... like listeners do

- **Hearing is *ecologically* grounded**
  - reflects 'natural scene' properties
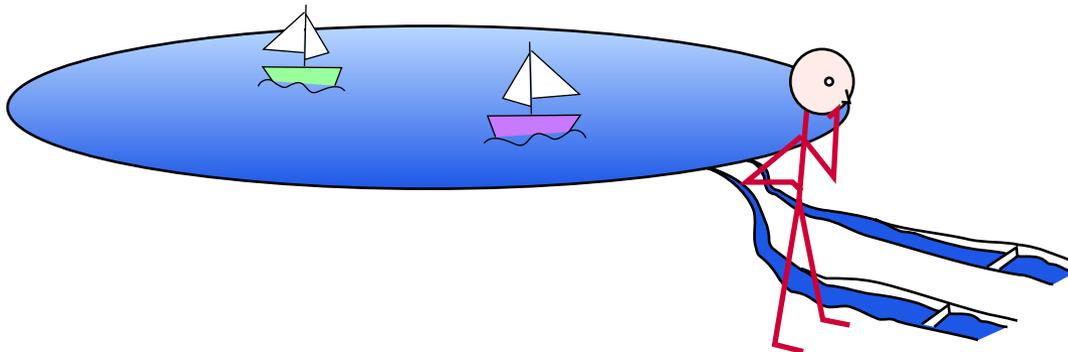  - subjective, not absolute

# Sound, mixtures, and learning



**Speech**   **Noise**   **Speech + Noise**

- **Sound**
  - carries useful information about the world
  - complements vision

- **Mixtures**
  - .. are the rule, not the exception
  - medium is 'transparent', sources are many
  - must be handled!

- **Learning**
  - the 'speech recognition' lesson:
    let the data do the work
  - like listeners

# The problem with recognizing 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)
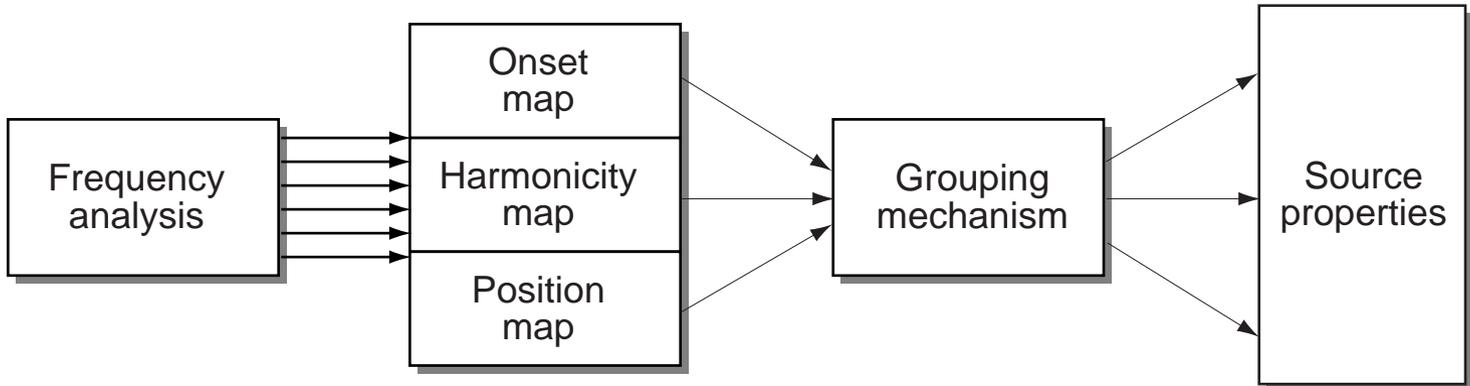
- **Received waveform is a mixture**
  - two sensors, N signals ... *underconstrained*

- **Disentangling mixtures as the primary goal?**
  - perfect solution is not possible
  - need experience-based *constraints*

Lab
ROSA
Laboratory for the Recognition and

# Human Auditory Scene Analysis
## (Bregman 1990)

- **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, ...):**
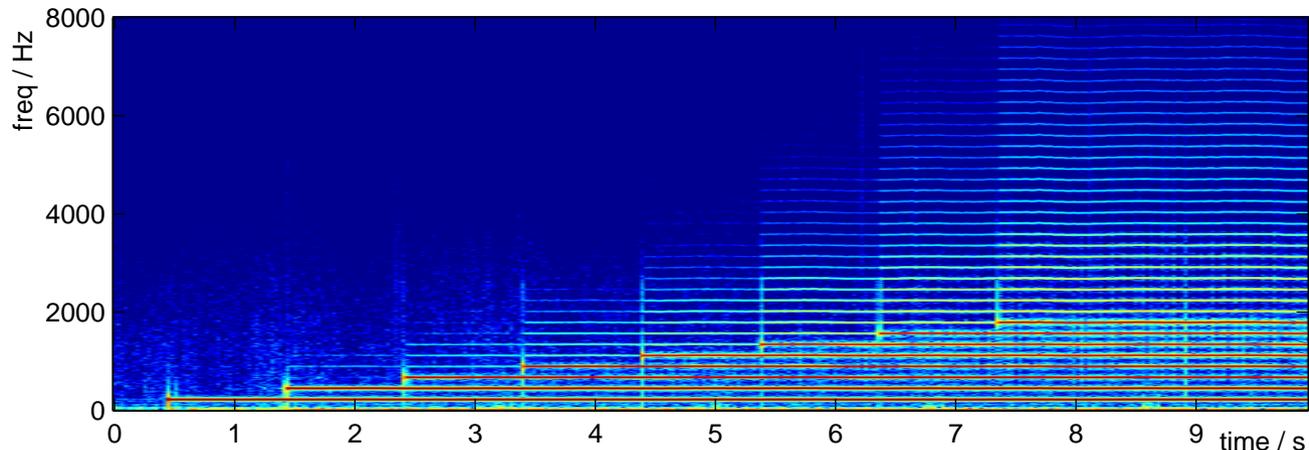  - cues: common onset/offset/modulation, harmonicity, spatial location, ...



*(after Darwin, 1996)*

Lab
ROSA
Laboratory for the Recognition and

# Cues to simultaneous grouping
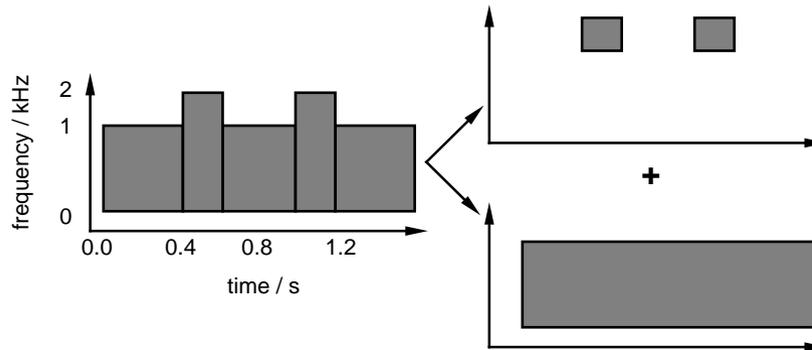
- **Elements + attributes**



- **Common onset**
  - simultaneous energy has common source

- **Periodicity**
  - energy in different bands with same cycle

- **Other cues**
  - spatial (ITD/IID), familiarity, ...

Lab
ROSA
Laboratory for the Recognition and
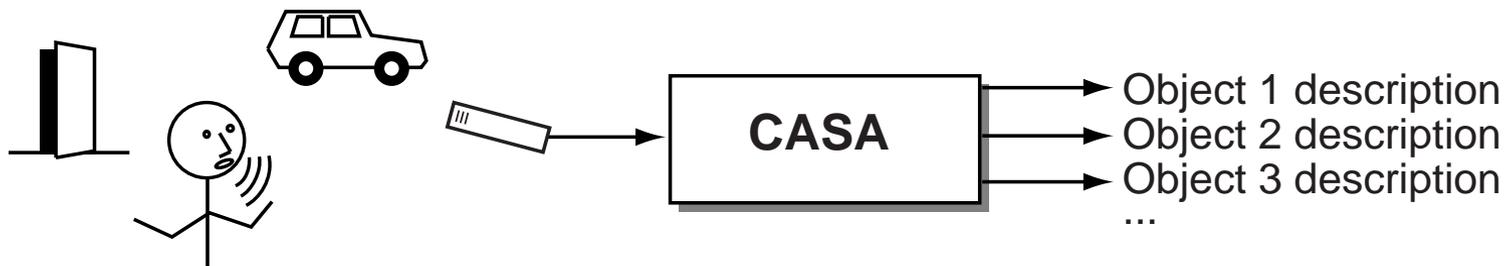
# The effect of context

- **Context can create an 'expectation':
  i.e. a bias towards a particular interpretation**

- **e.g. Bregman's "old-plus-new" principle:**

  A change in a signal will be interpreted as an
  *added* source whenever possible



- a different division of the same energy
  depending on what preceded it

Lab
ROSA
Laboratory for the Recognition and

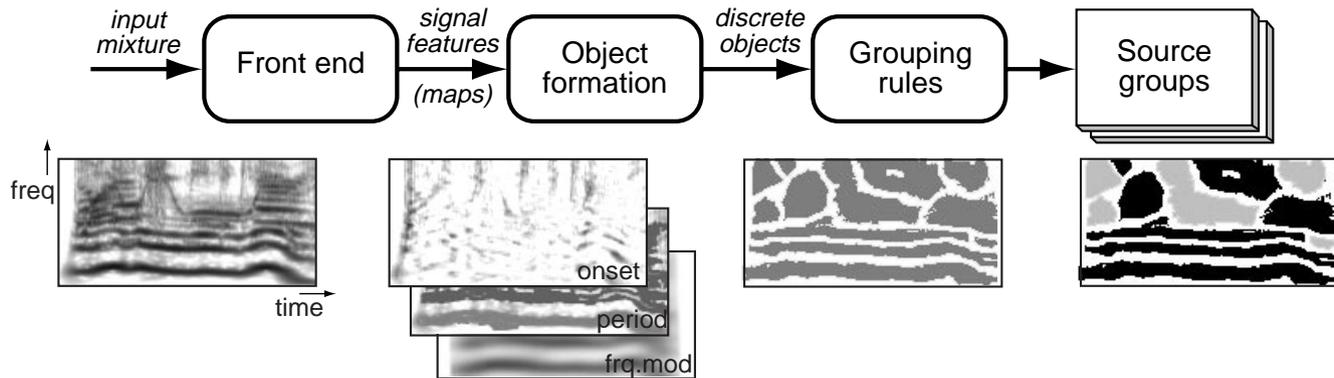# Computational Auditory Scene Analysis (CASA)



- **Goal: Automatic sound organization ;
  Systems to 'pick out' sounds in a mixture**
  - ... like people do

- **E.g. voice against a noisy background**
  - to improve speech recognition

- **Approach:**
  - psychoacoustics describes grouping 'rules'
  - ... just implement them?

Lab
ROSA
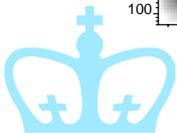Laboratory for the Recognition and

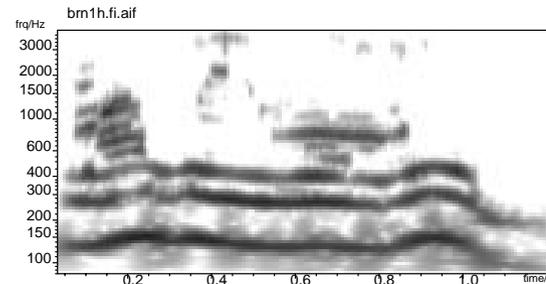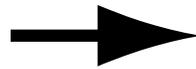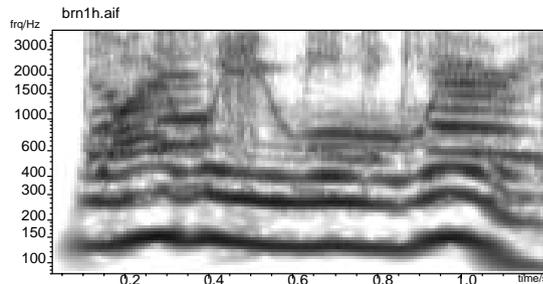# The Representational Approach
## (Brown & Cooke 1993)
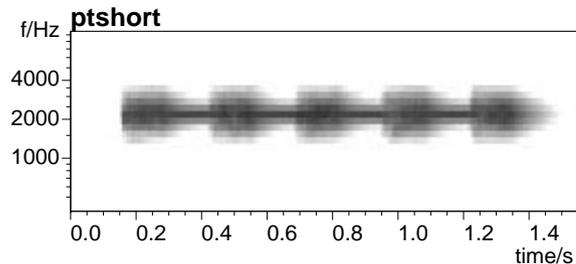
- **Implement psychoacoustic theory**



- 'bottom-up' processing
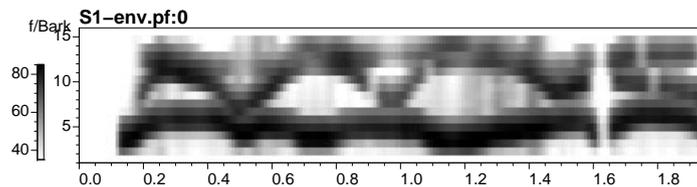- uses common onset & periodicity cues

- **Able to extract voiced speech:**

# Restoration in sound perception

- **Auditory 'illusions' = hearing what's not there**

- **The continuity illusion**

**ptshort**

f/Hz

4000
2000
1000

0.0   0.2   0.4   0.6   0.8   1.0   1.2   1.4
time/s

- **SWS**

**S1−env.pf:0**

f/Bark
15
80
10
60
5
40

0.0  0.2  0.4  0.6  0.8  1.0  1.2  1.4  1.6  1.8

- duplex perception
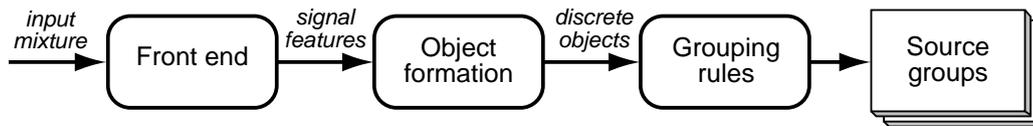
- **How to model in CASA?**

Lab
ROSA
Laboratory for the Recognition and
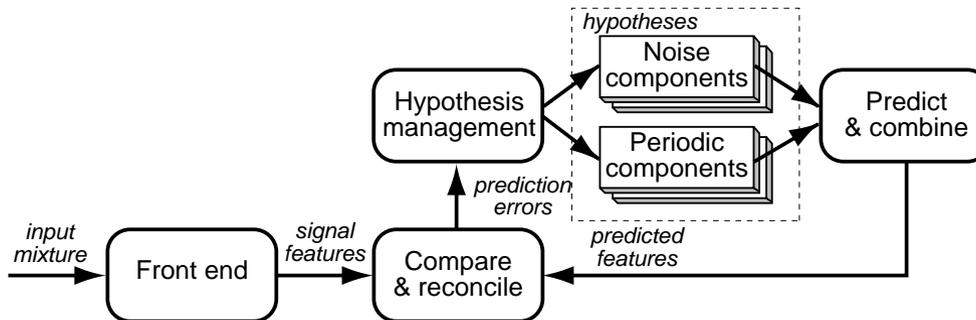
# Adding top-down constraints

**Perception is not *direct*
but a *search* for *plausible hypotheses***

- **Data-driven (bottom-up)...**



- objects irresistibly appear
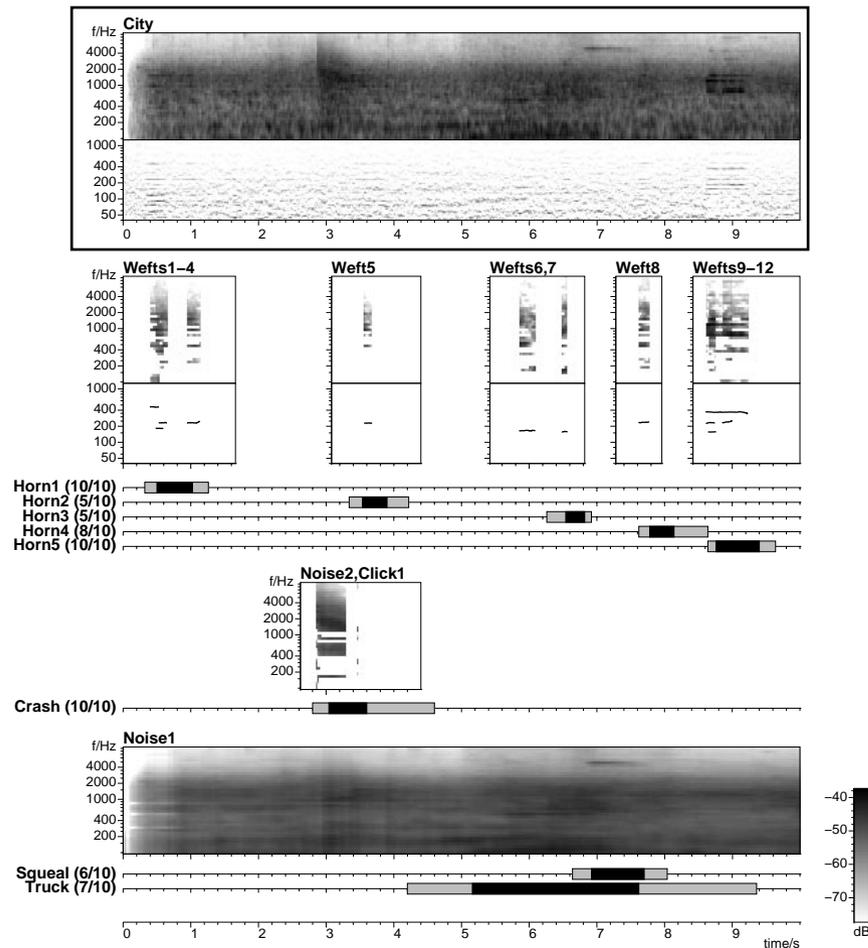
## vs. Prediction-driven (top-down)



- match observations
  with parameters of a world-model
- need world-model constraints...

Lab
ROSA
Laboratory for the Recognition and

# Prediction-Driven CASA
## (Ellis 1996)

- **Explain a complex sound with basic elements**

# Approaches to sound mixture recognition

- **Recognize combined signal**
  - 'multicondition training'
  - combinatorics..

- **Separate signals**
  - e.g. CASA, ICA
  - nice, if you can do it

- **Segregate features into fragments**
  - then missing-data recognition

Lab
ROSA
Laboratory for the Recognition and

# Aside: Evaluation

- **Evaluation is a big problem for CASA**
  - what is the goal, really?
  - what is a good test domain?
  - how do you measure performance?

- **SNR improvement**
  - not easy given only before-after signals: correspondence problem
  - can do with fixed filtering mask; rewards removing signal as well as noise

- **ASR improvement**
  - recognizers typically very sensitive to artefacts

- **'Real' task?**
  - mixture corpus with specific sound events...

Lab
ROSA
Laboratory for the Recognition and

# Outline

**1** Auditory Scene Analysis

**2** **Speech Recognition & Mixtures**
- - the information in speech
- - Meeting Recorder project
- - speech fragment decoding

**3** Music Analysis & Similarity

**4** General Sound Organization

**5** Future Work

Lab
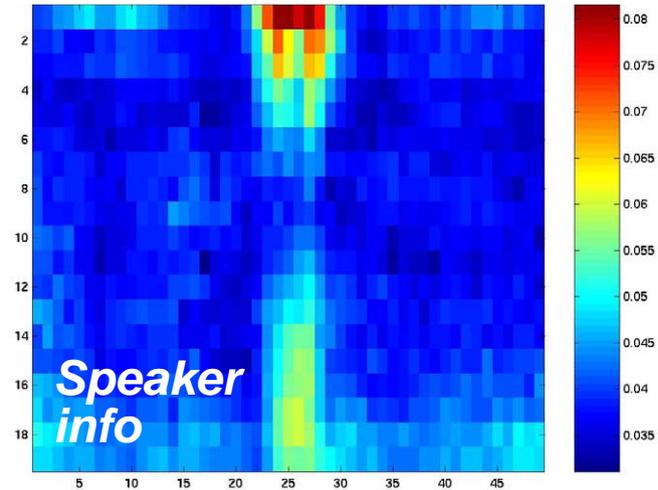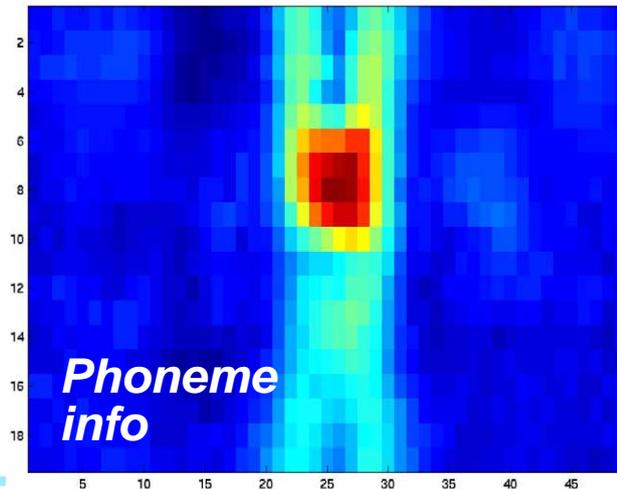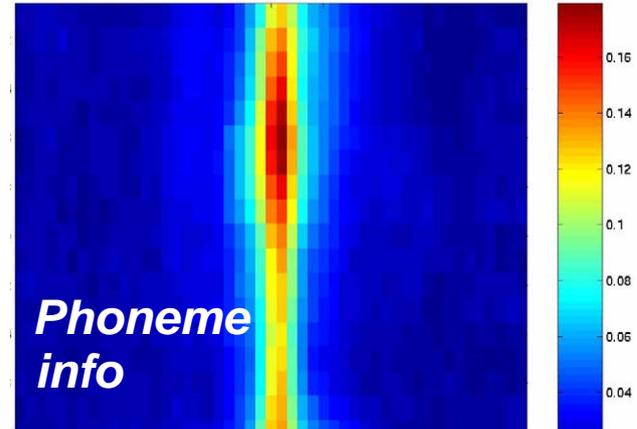ROSA
Laboratory for the Recognition and

# The information in speech

### (Patricia Scanlon)
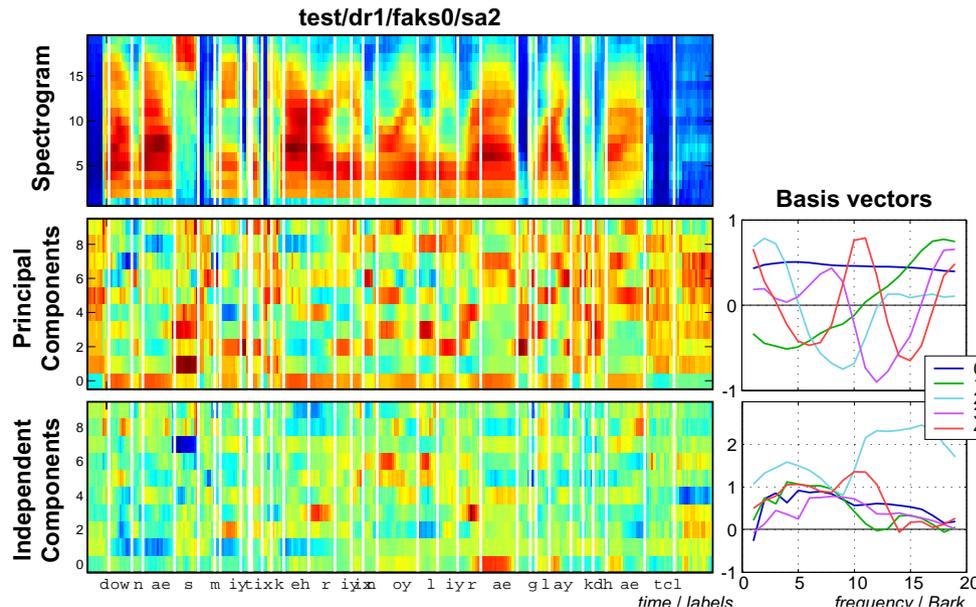
- **Mutual Information identifies where the information is in time/ frequency:**
  - little temporal structure averaged over all sounds

  *Phoneme info*

  - **Better with just *vowels:***

  *Phoneme info*

  *Speaker info*

LAB
ROSA
Laboratory for the Recognition and

# The best subword units?
### (Eric Fosler)

- **Speech recognizers typically use phonemes**
  - inherited from linguistics

- **Alternative approach is 'articulatory features'**
  - orthogonal attributes defining subwords

- **Can we infer a feature set from the data**
  - using e.g. Independent Component Analysis

# The Meeting Recorder Project
(CompSci, ICSI, UW, IDIAP, SRI, IBM)

- **Microphones in conventional meetings**
  - for summarization/retrieval/behavior analysis
  - informal, overlapped speech

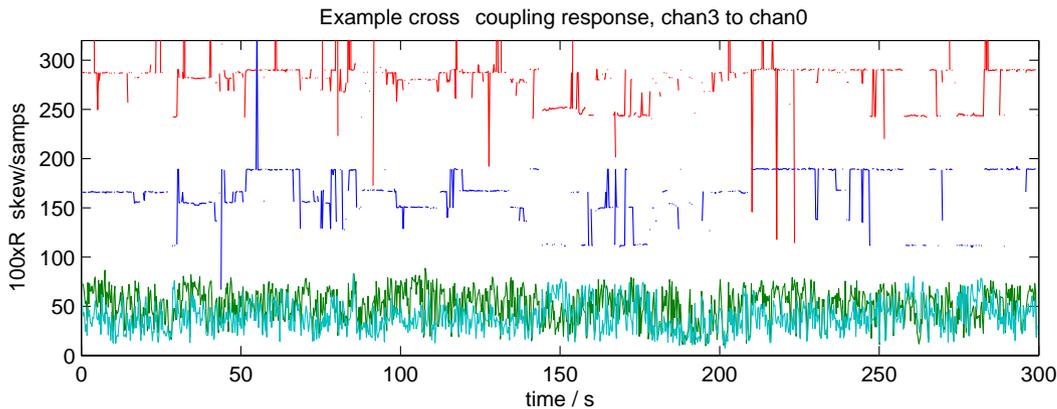- **Data collection (ICSI, UW, IDIAP):**



  - 100 hours collected, ongoing transcription

- **NSF 'Mapping Meetings' project**
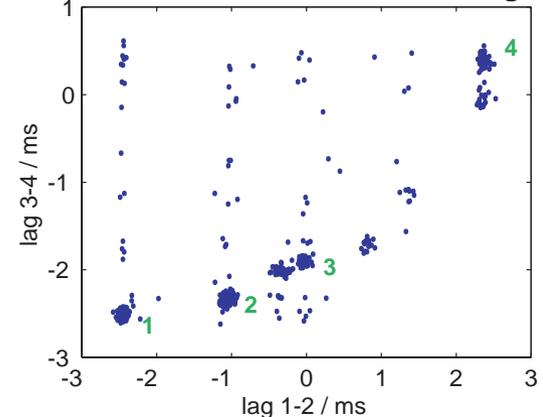  - also interest from NIST, DARPA, EU

Lab
ROSA
Laboratory for the Recognition and

# Speaker Turn detection
## (Huan Wei Hee, Jerry Liu)

- **Acoustic:**
  **Triangulate tabletop mic timing differences**
  - use normalized peak value for confidence



Example cross coupling response, chan3 to chan0



mr-2000-11-02-1440: PZM xcorr lags

- **Behavioral: Look for patterns of speaker turns**



mr04: Hand-marked speaker turns vs. time + auto/manual boundaries

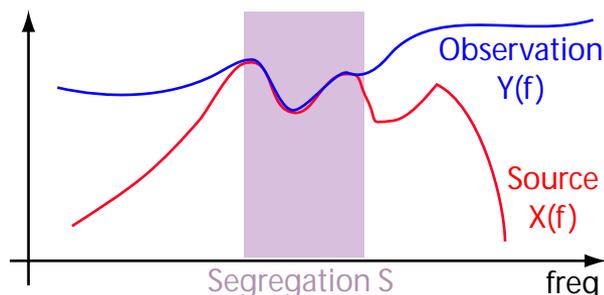Lab ROSA
Laboratory for the Recognition and

# Speech Fragment recognition
## (Barker & Cooke/Sheffield)

- **Standard classification chooses between models $M$ to match source features $X$**

$$M^* = \underset{M}{\mathrm{argmax}} \, P(M|X) = \underset{M}{\mathrm{argmax}} \, P(X|M) \cdot \frac{P(M)}{\cancel{P(X)}}$$

- **Mixtures → observed features $Y$, segregation $S$, all related by $P(X|Y, S)$**



Observation Y(f)

Source X(f)

Segregation S          freq
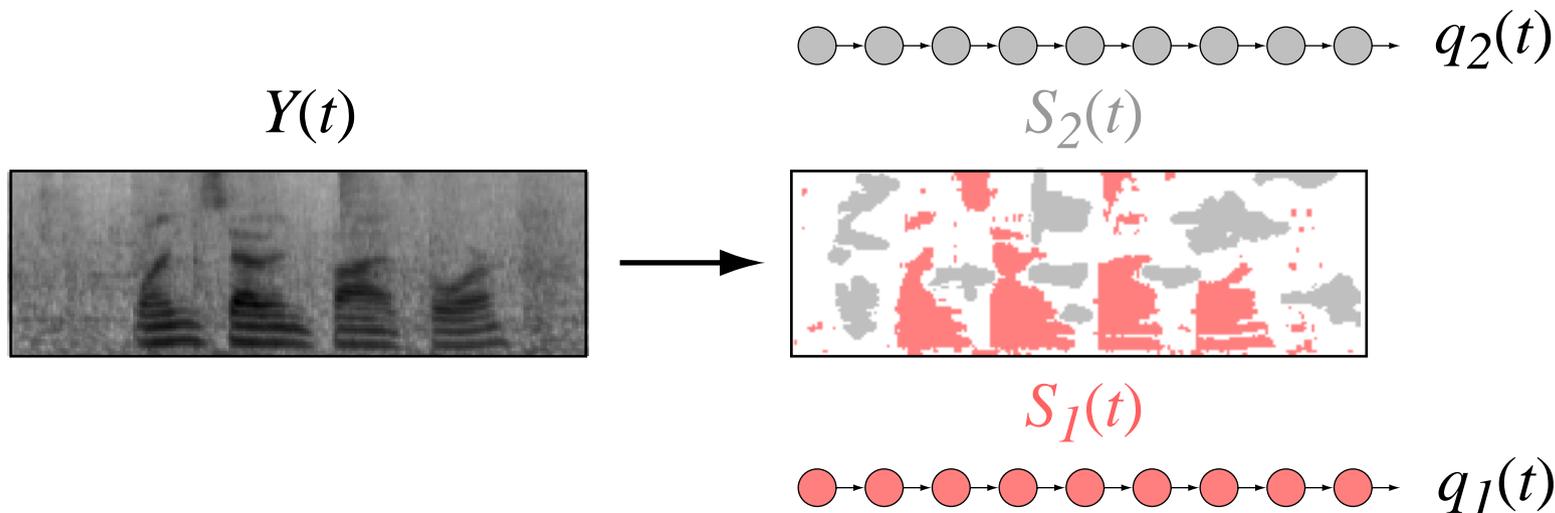
- *spectral features* allow clean relationship

- **Joint classification of model and segregation:**

$$P(M, S|Y) = P(M)\int P(X|M) \cdot \frac{P(X|Y, S)}{P(X)} dX \cdot P(S|Y)$$

# Multi-source decoding

- **Search for more than one source**



$q_2(t)$

$Y(t)$

$S_2(t)$

$S_1(t)$

$q_1(t)$

- **Mutually-dependent data masks**

- **Use e.g. CASA features to propose masks**
  - locally coherent regions

- **Theoretical vs. practical limits**

Lab
ROSA
Laboratory for the Recognition and

# Outline

**1** **Auditory Scene Analysis**

**2** **Speech Recognition & Mixtures**

**3** **Music Analysis & Similarity**
- musical structure analysis
- similarity browsing

**4** **General Sound Organization**

**5** **Future Work**
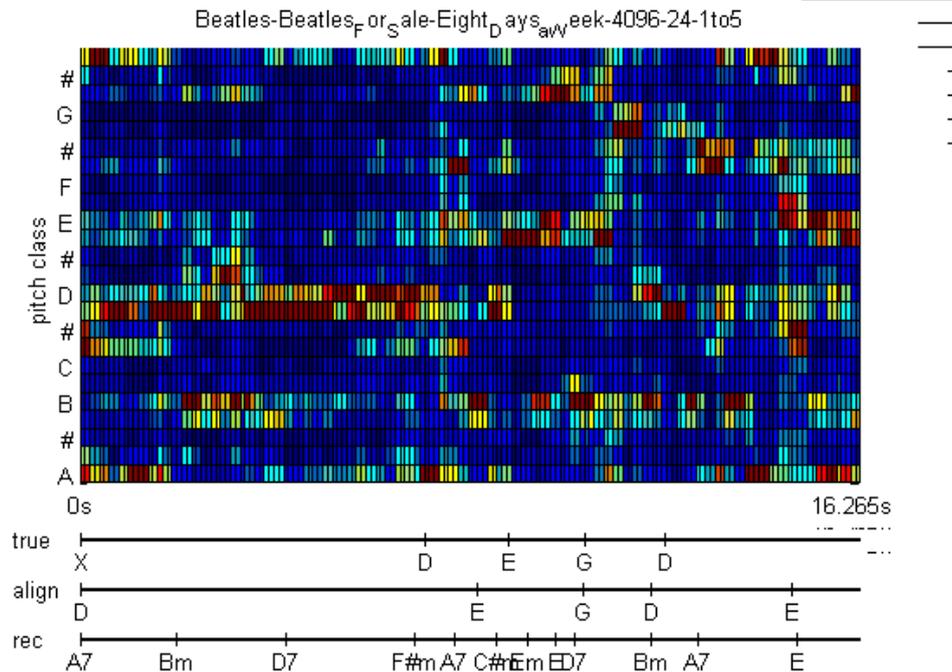
Lab
**ROSA**
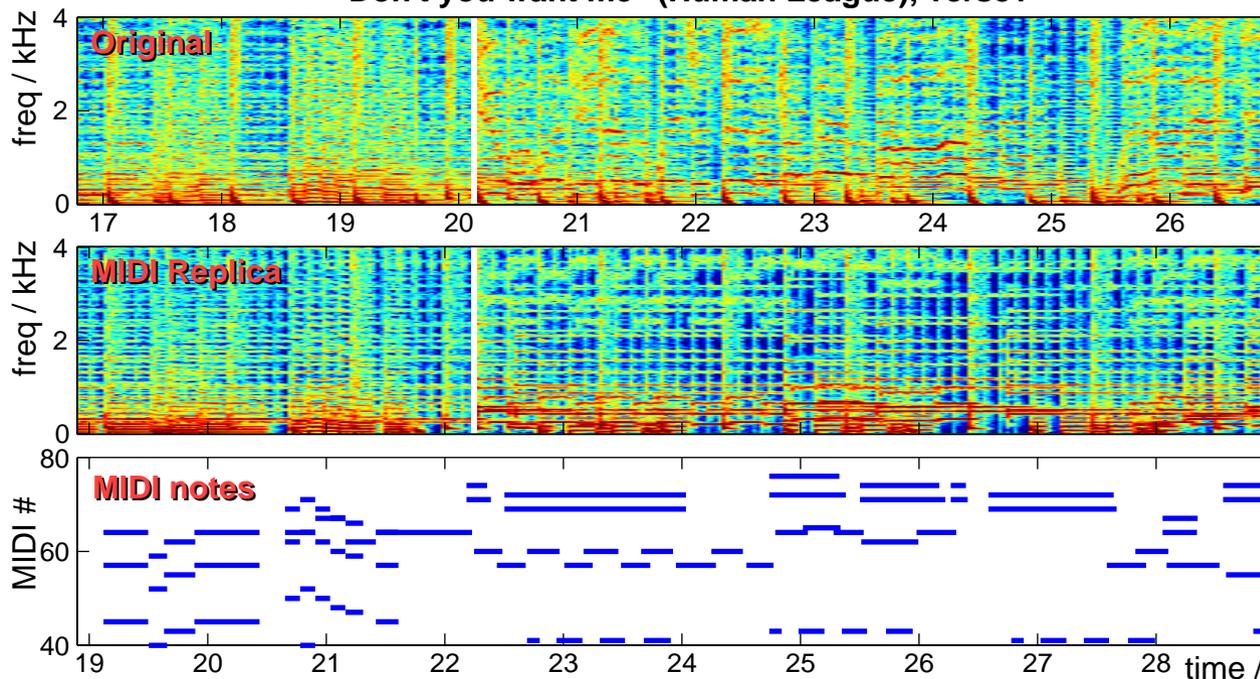Laboratory for the Recognition and

# Music Structure Analysis

## (Alex Sheh)

- **Fine-level information from music**
  - for searching
  - for modeling/statistics

- **e.g. Chord sequences via PCPs :**



Beatles-Beatles$_F$or$_S$ale-Eight$_D$ays$_{a}$w$_{eek}$-4096-24-1to5

# Ground truth for Music Recordings

(Rob Turetsky)

- **Machine Learning algorithms need labels**
  - but real recordings don't have labels

- **MIDI 'replicas' exist**

- **Alignment locates MIDI notes in real sound:**

"Don't you want me" (Human League), verse1

LAB
ROSA
Laboratory for the Recognition and

# Music Similarity Browsing
## (Adam Berenzweig)

- **'Anchor models' : music on subjective axes**

ROSA
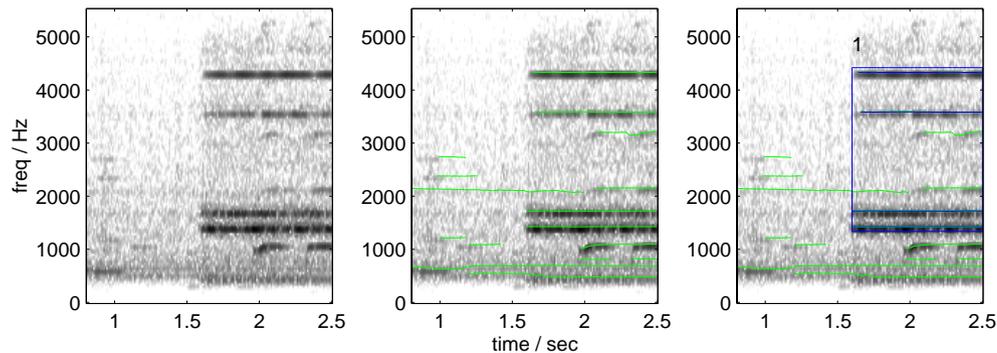Laboratory for the Recognition and

# Outline

**1** Auditory Scene Analysis

**2** Speech Recognition & Mixtures

**3** Music Analysis & Similarity

**4** General Sound Organization
- alarm detection
- sound texture modeling
- recognition of multiple sources

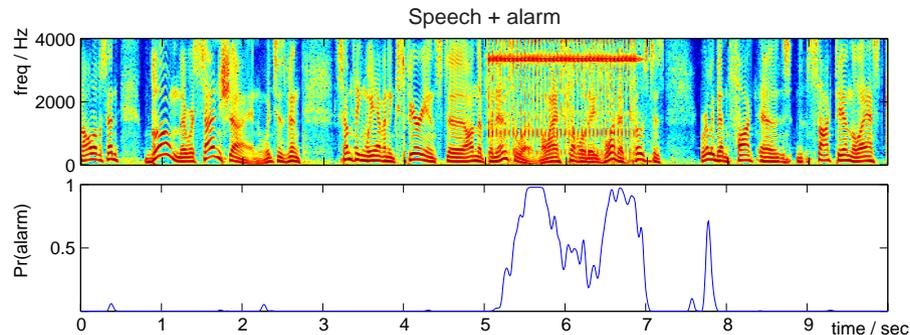**5** Future Work

Lab
ROSA
Laboratory for the Recognition and

# Alarm sound detection

- **Alarm sounds have particular structure**
  - people 'know them when they hear them'

- **Isolate alarms in sound mixtures**
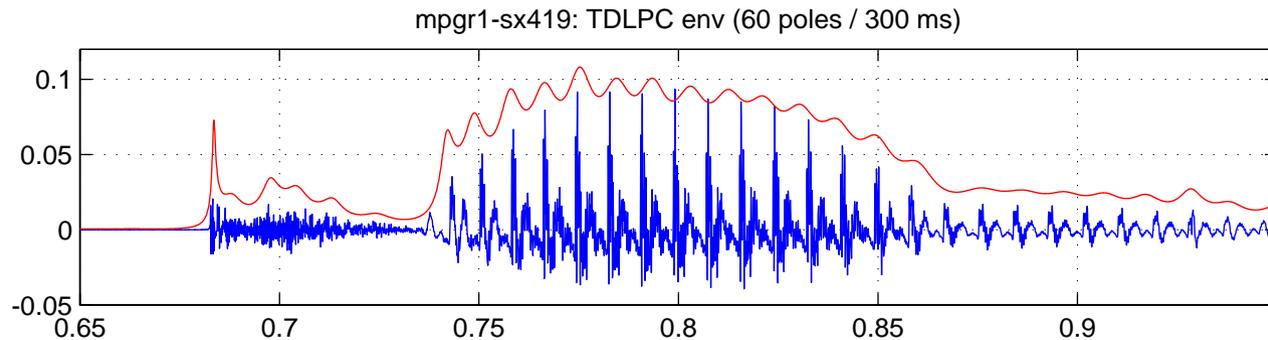


- sinusoid peaks have invariant properties

Speech + alarm



- cepstral coefficients are easy to model

Lab
ROSA
Laboratory for the Recognition and

# Sound Texture Modeling

(Marios Athineos)

- **Best sound models are based on sinusoids**
  - noise residual modeled quite simply

- **Noise 'textures' have extra temporal structure**
  - need a more detailed model

- **Linear prediction of spectrum defines a parametric temporal envelope:**

mpgr1-sx419: TDLPC env (60 poles / 300 ms)



- **High-quality noise-excited resynthesis:**
  - original  -  resynth  -  x2 TSM  -  c/w PVOC

LabROSA
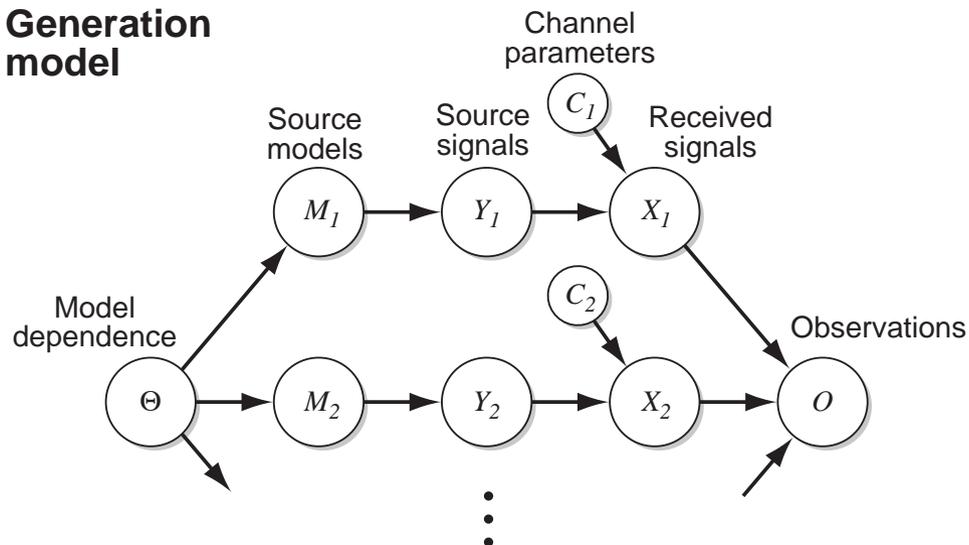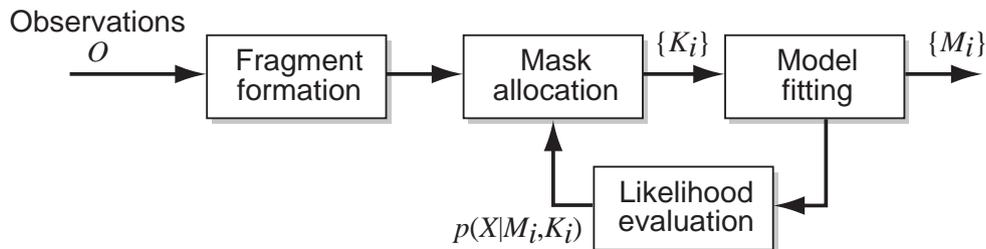Laboratory for the Recognition and

# Sound mixture decomposition

## (Manuel Reyes)

- **Full or approximate Bayesian inference to model multiple, independent sound sources:**

**Generation model**



**Analysis structure**

Lab
ROSA
Laboratory for the Recognition and

# Outline

**1** **Auditory Scene Analysis**

**2** **Speech Recognition & Mixtures**

**3** **Music Analysis & Similarity**

**4** **General Sound Organization**

**5** **Future Work**
- audio-visual information
- real-world sound indexing

Lab
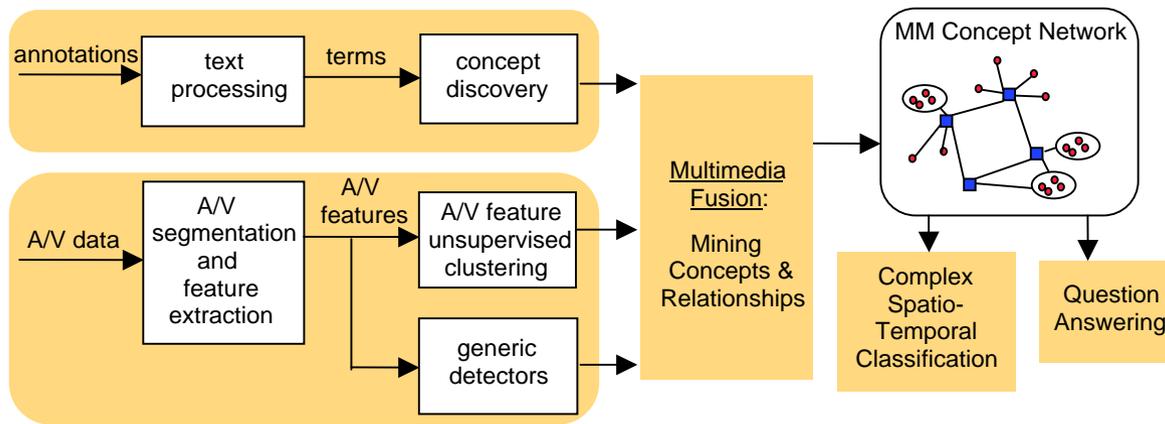ROSA
Laboratory for the Recognition and

# Future work:
# Automatic audio-video analysis
### (Shih-Fu Chang, Kathy McKeown)

- **Documentary archive management**
  - huge ratio of raw-to-finished material
  - costly manual logging

- **Problem: term ↔ signal mapping**
  - training corpus of past annotations
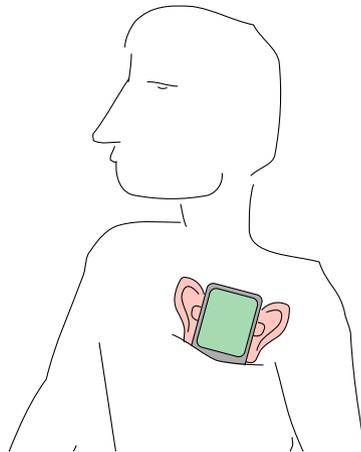  - interactive semi-automatic learning

# The 'Listening Machine'

- **Smart PDA records everything**

- **Only useful if we have index, summaries**
  - monitor for particular sounds
  - real-time description

- **Scenarios**

  - personal listener $\rightarrow$ summary of your day
  - future prosthetic hearing device
  - autonomous robots

- **Meeting data, ambulatory audio**

# LabROSA Summary

**DOMAINS**

- Broadcast
- Movies
- Lectures

- Meetings
- Personal recordings
- Location monitoring

## ROSA

- Object-based structure discovery & learning

- Speech recognition
- Speech characterization
- Nonspeech recognition

- Scene analysis
- Audio-visual integration
- Music analysis

**APPLICATIONS**

- Structuring
- Search
- Summarization
- Awareness
- Understanding

Lab
ROSA
Laboratory for the Recognition and