

---

---

# Recognition and Organization of Speech and Audio

Dan Ellis  
<dpwe@ee.columbia.edu>

Laboratory for Recognition and Organization of Speech and Audio  
(Lab**ROSA**)

Electrical Engineering, Columbia University  
<http://labrosa.ee.columbia.edu/>

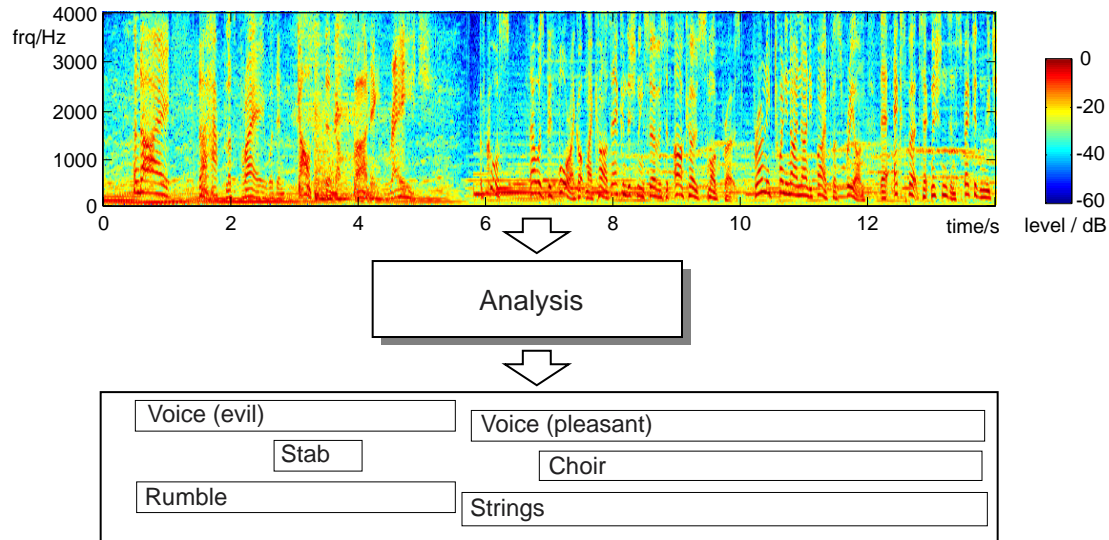
## Outline

- 1 Audio Organization
- 2 Speech, music, and other
- 3 Future work



## 1

# Audio Organization

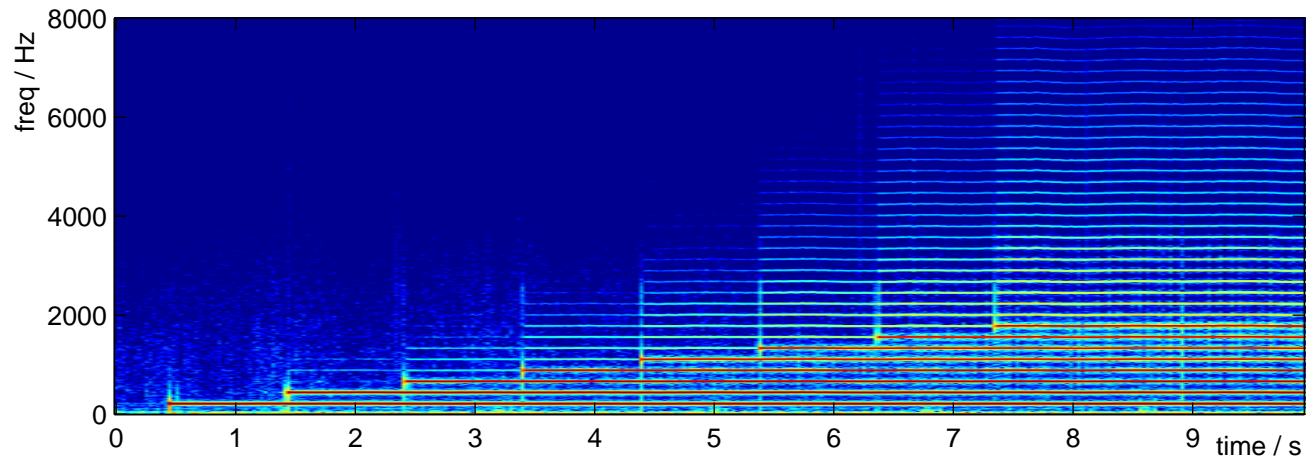


- **Analyzing and describing complex sounds:**
  - continuous sound mixture  
→ distinct objects & events
- **Human listeners as the prototype**
  - strong subjective impression when listening
  - ..but hard to 'see' in signal



# Human Sound Organization

- **Sound percepts depend on ‘grouping cues’**
  - common onset across frequency
  - common periodicity (fundamental)
  - spatial (binaural) cues, familiarity, ...



- **Hearing confers evolutionary advantage**
  - optimized to get ‘useful’ information from sound
- **Auditory perception is *ecologically* grounded**
  - scene analysis is preconscious (→ illusions)



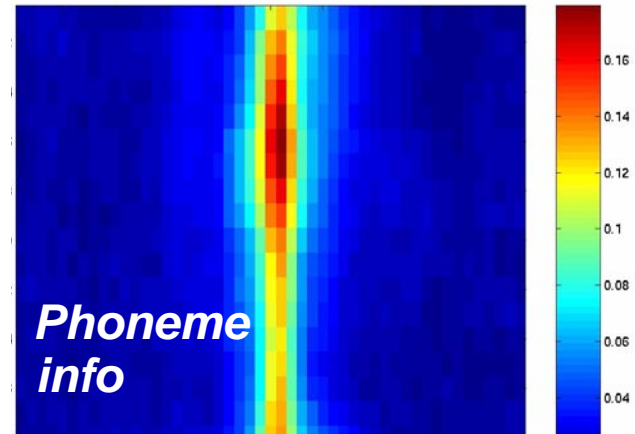
## 2

# Speech & Speaker recognition

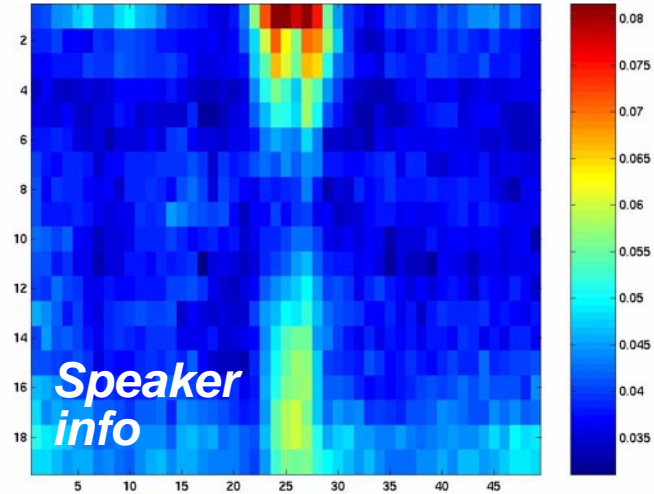
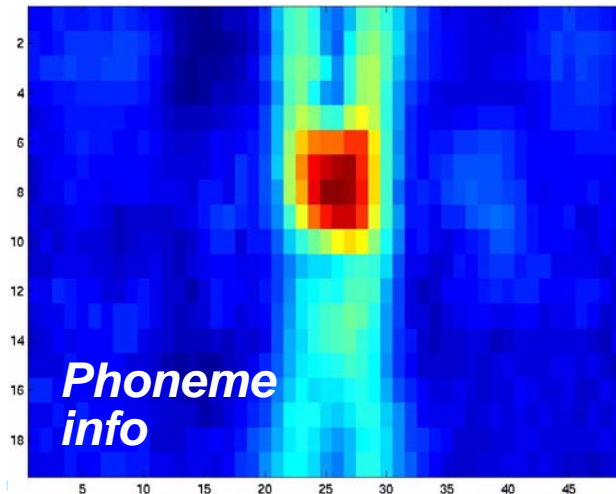
(Patricia Scanlon)

- **Mutual Information identifies where the information is in time/frequency:**

- little temporal structure averaged over all sounds



- **Better with just vowels:**



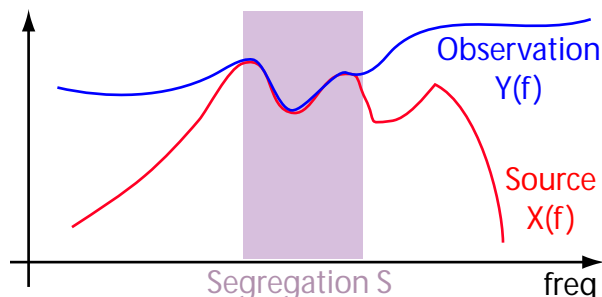
# Speech Fragment recognition

(Barker & Cooke/Sheffield)

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

$$M^* = \operatorname{argmax}_M P(M|X) = \operatorname{argmax}_M P(X|M) \cdot \frac{P(M)}{P(X)}$$

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



- *spectral features* allow clean relationship

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

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



---

---

# 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):**



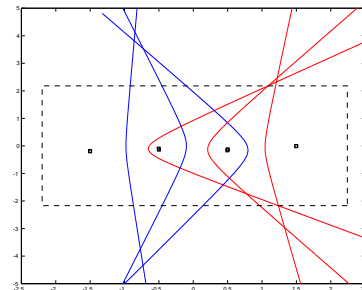
- 100 hours collected, ongoing transcription
- **NSF ‘Mapping Meetings’ project**
  - also interest from NIST, DARPA, EU



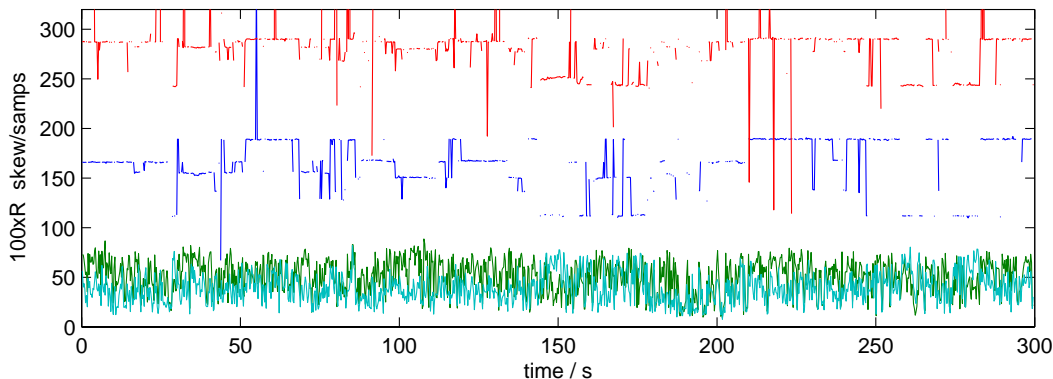
# Tabletop mics: Turn detection

(Huan Wei Hee, Jerry Liu)

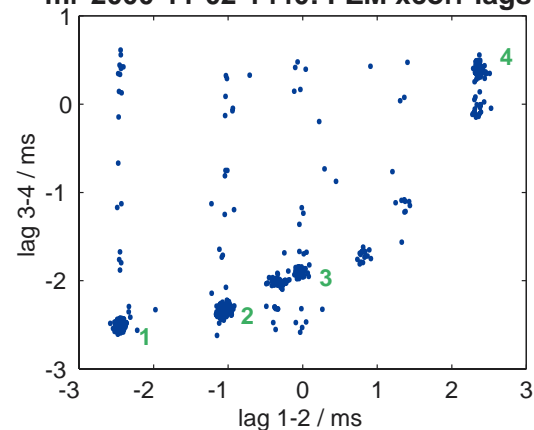
- **4 mics ~ 1m separated along center of table**
  - 3 timing differences
  - slight U/D offset to disambiguate
- **Hi-res cross-correlation for timings**
  - use normalized peak value for confidence
  - cluster results



Example cross coupling response, chan3 to chan0



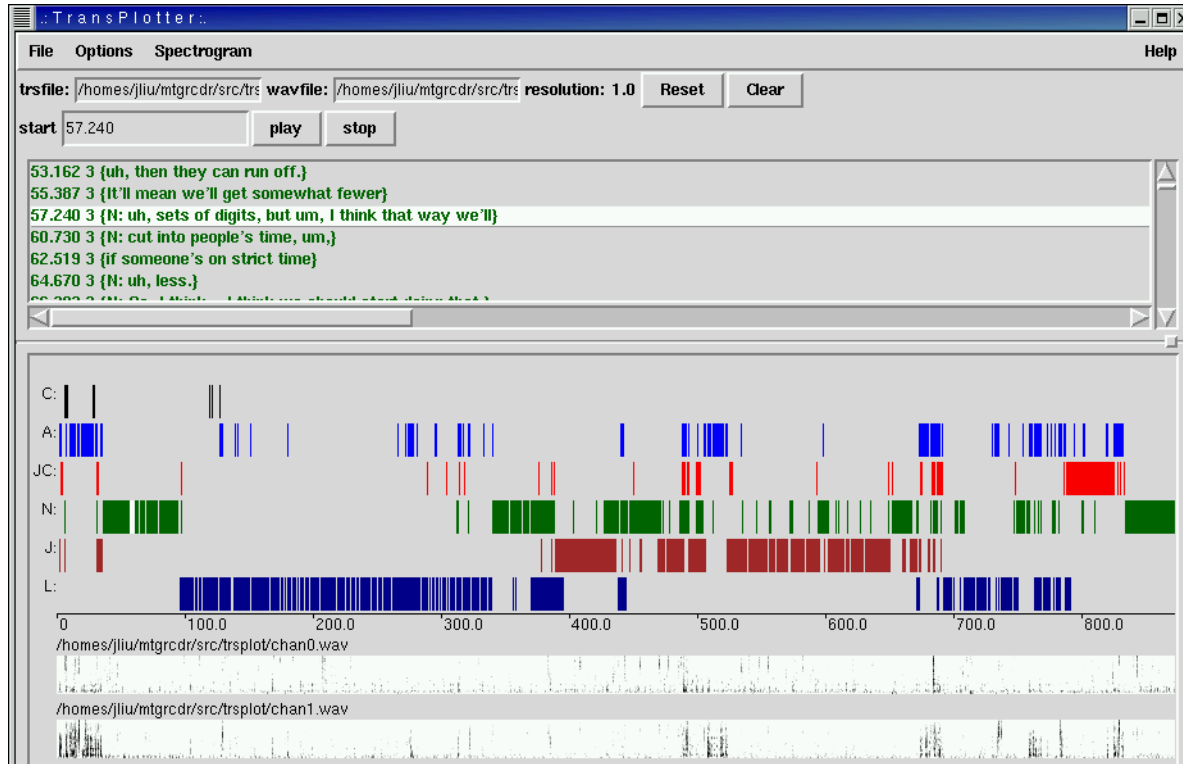
mr-2000-11-02-1440: PZM xcorr lags



# Visualization: transPlotter

(Jerry Liu)

- Speaker turn *patterns* are informative



- **Browser for 'high-level' view, quick examination**
  - snack, iwidgets based

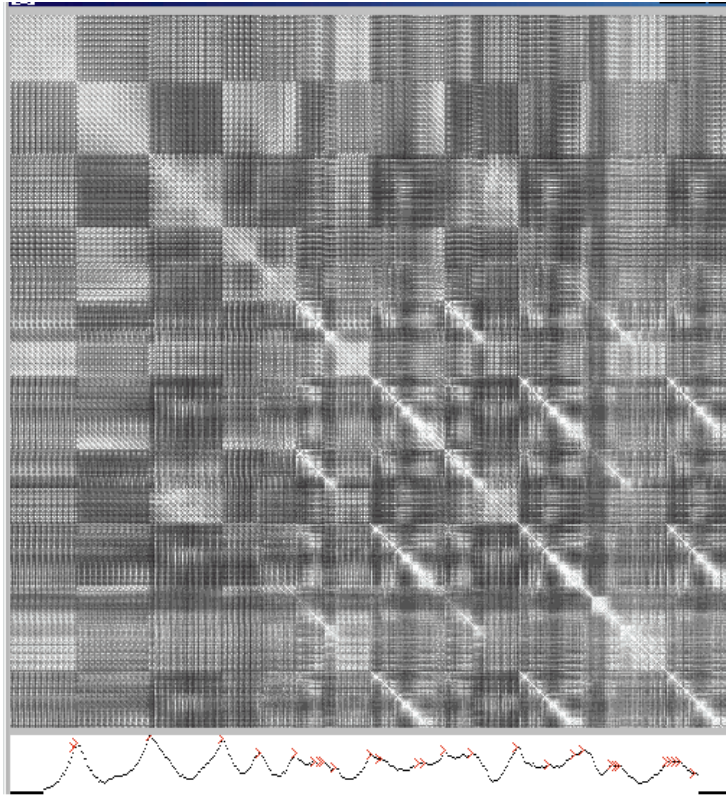




# Music analysis: Structure recovery

(Rob Turetsky, Alex Sheh)

- **Identify & match segment-level structure in music**



- Foote similarity matrices point to repeated sections
- Chord segmentation using ASR 'hidden state' model (train with chord transcriptions)
- Note transcription/timing by aligning audio to resynthesized MIDI versions



# Music Similarity & Recommendation

(Adam Berenzweig, Lawrence/NECI)

**Playola** Search:  Artist   
[\[About\]](#) [\[Help\]](#) [\[Turn Samples Off\]](#) [\[Turn Debug On\]](#) [\[Turn Poppers Off\]](#) [\[Logout dpwe\]](#)

Get Playola Selections: 20 songs  you recently heard   Browse: [Artists](#) [Albums](#) [Playlists](#) Range: 0-C

Artist: **The Woodbury Muffin Outbreak** [\[band web page\]](#) [\[Play!\]](#) Playlist: -New Playlist-  [\[Add to\]](#) [\[View\]](#)

	Song Title	Artist	Time	Rating
<input type="checkbox"/>	The Ballad of Tabitha	<a href="#">The Woodbury Muffin Outbreak</a>	4:00	<input type="checkbox"/>
<input type="checkbox"/>	Monkey Dreams	<a href="#">The Woodbury Muffin Outbreak</a>	2:57	<input type="checkbox"/>
<input type="checkbox"/>	A Cold Dark Night (Live)	<a href="#">The Woodbury Muffin Outbreak</a>	3:13	<input type="checkbox"/>
<input type="checkbox"/>	Leo, The Ballad of	<a href="#">The Woodbury Muffin Outbreak</a>	1:48	<input type="checkbox"/>
<input type="checkbox"/>	Baby I Forgot To Tell You	<a href="#">The Woodbury Muffin Outbreak</a>	4:04	<input type="checkbox"/>

**Music-Space Browser** [\[What's This?\]](#)

Feature	Less	More
AltNGrunge	<input type="checkbox"/>	<input type="checkbox"/>
CollegeRock	<input type="checkbox"/>	<input type="checkbox"/>
Country	<input type="checkbox"/>	<input type="checkbox"/>
DanceRock	<input type="checkbox"/>	<input type="checkbox"/>
Electronica	<input type="checkbox"/>	<input type="checkbox"/>
MetalNPunk	<input type="checkbox"/>	<input type="checkbox"/>
NewWave	<input type="checkbox"/>	<input type="checkbox"/>
Rap	<input type="checkbox"/>	<input type="checkbox"/>
RnBSoul	<input type="checkbox"/>	<input type="checkbox"/>
SingerSongwriter	<input type="checkbox"/>	<input type="checkbox"/>
SoftRock	<input type="checkbox"/>	<input type="checkbox"/>
TradRock	<input type="checkbox"/>	<input type="checkbox"/>
Female	<input type="checkbox"/>	<input type="checkbox"/>
Hifi	<input type="checkbox"/>	<input type="checkbox"/>

**Similar Songs:** [\[Play this list\]](#) [\[What's This?\]](#)

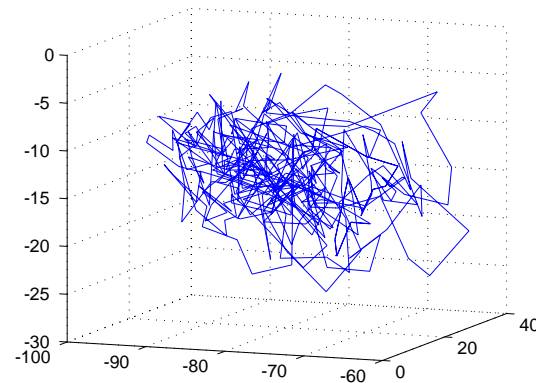
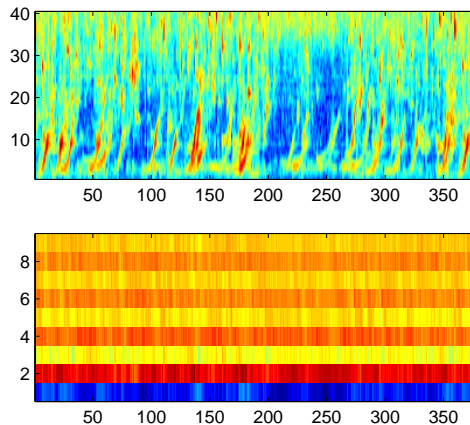
	Song Title	Artist	Distance	Good Match?
<input type="checkbox"/>	Baby I Forgot To Tell You	<a href="#">The Woodbury Muffin Outbreak</a>	0.00	<input type="checkbox"/>
<input type="checkbox"/>	Number five	<a href="#">Bizi Chyld</a>	0.07	<input type="checkbox"/>
<input type="checkbox"/>	Waiting for Your Love	<a href="#">Toto</a>	0.08	<input type="checkbox"/>
<input type="checkbox"/>	Except from ICB	<a href="#">Weirdmusic</a>		<input type="checkbox"/>



# Sound texture modeling

(Marios Athineos)

- **Sound textures are important but neglected**
  - fire, rain, paper: no clear pitch, onsets, shape
  - typically modeled with filtered white noise
- **LPC on *spectrum* captures temporal structure:**
  - better parameterization of texture structure?



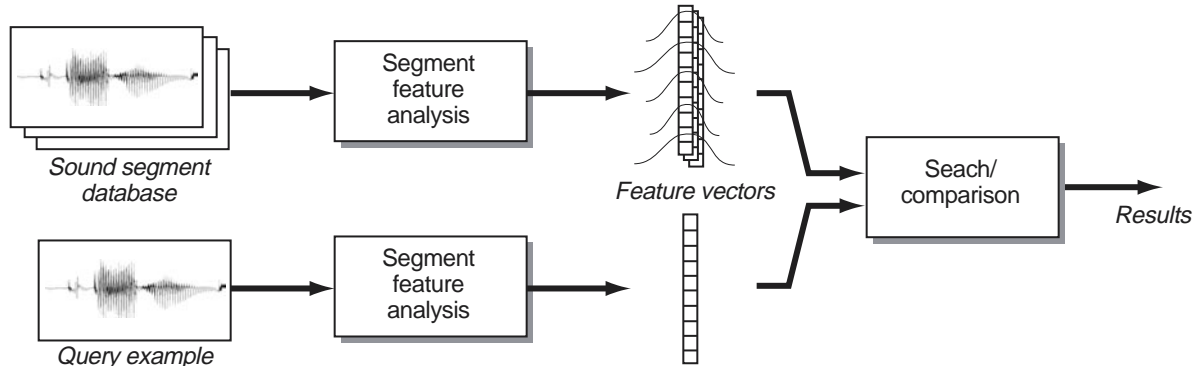
- generate variations?



# Audio Information Retrieval

(Manuel Reyes)

- **Searching in a database of audio**
  - speech .. use ASR
  - text annotations .. search them
  - sound effects library?
- **e.g. Muscle Fish “SoundFisher” browser**
  - define multiple ‘perceptual’ feature dimensions
  - search by proximity in (weighted) feature space

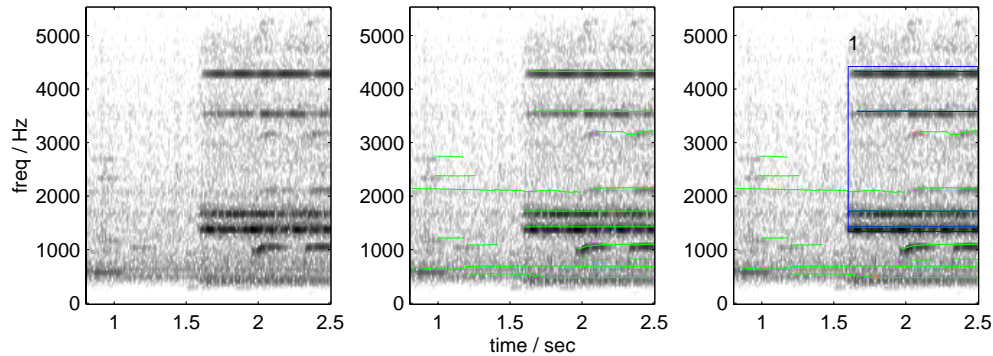


- features are ‘global’ for each soundfile,  
no attempt to separate mixtures

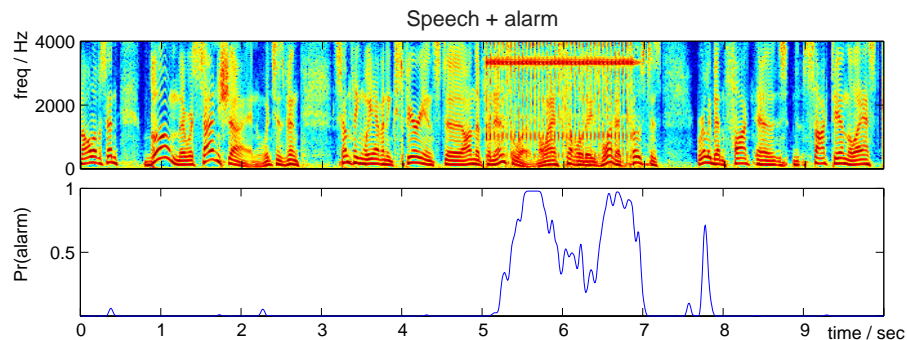


# Alarm sound detection

- **Alarm sounds have particular structure**
  - people 'know them when they hear them'
- **Isolate alarms in sound mixtures**



- sinusoid peaks have invariant properties



- cepstral coefficients are easy to model

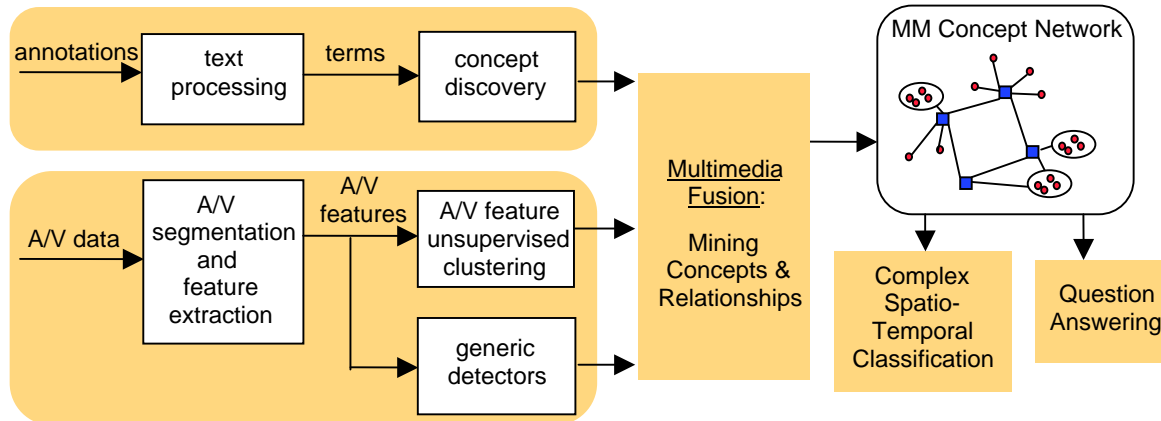


## 2

# 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



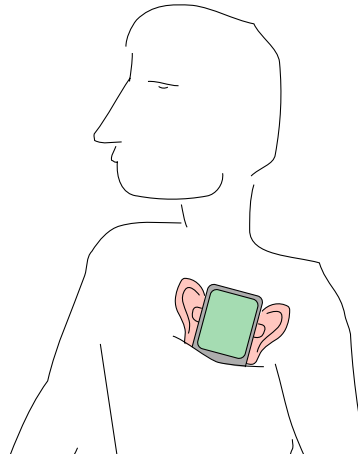


---

---

# The 'Machine listener'

- **Goal: An auditory system for machines**
  - use same environmental information as people
- **Signal understanding**
  - monitor for particular sounds
  - real-time description
- **Scenarios**



- personal listener → summary of your day
- future prosthetic hearing device
- autonomous robots



# LabROSA Summary

## DOMAINS

- Broadcast
- Meetings
- Movies
- Personal recordings
- Lectures
- Location monitoring

## ROSA

- Object-based structure discovery & learning
- Speech recognition
- Scene analysis
- Speech characterization
- Audio-visual integration
- Nonspeech recognition
- Music analysis

## APPLICATIONS

- Structuring
- Search
- Summarization
- Awareness
- Understanding

