# LabROSA
# Research Overview

Dan Ellis

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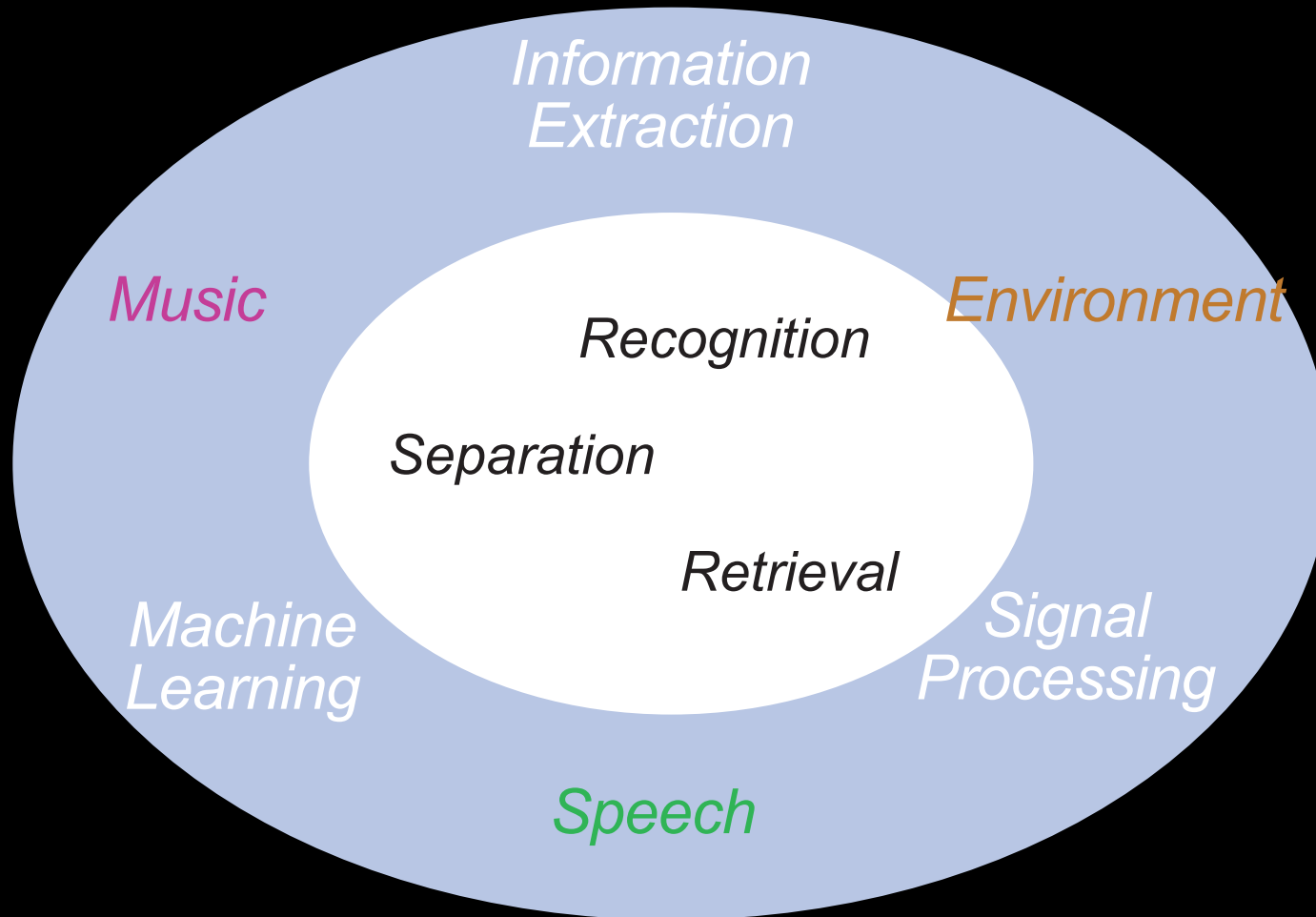
dpwe@ee.columbia.edu          http://labrosa.ee.columbia.edu/

1. Music
2. Environmental sound
3. Speech Enhancement

Lab
ROSA

Laboratory for the Recognition and
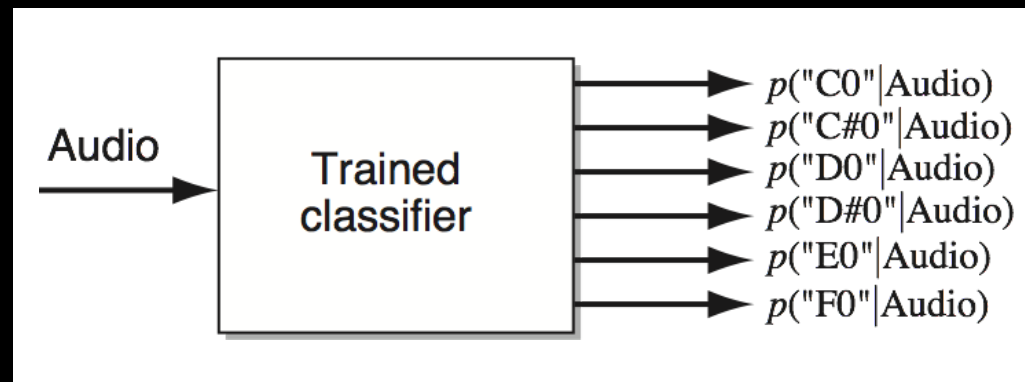Organization of Speech and Audio

COLUMBIA UNIVERSITY
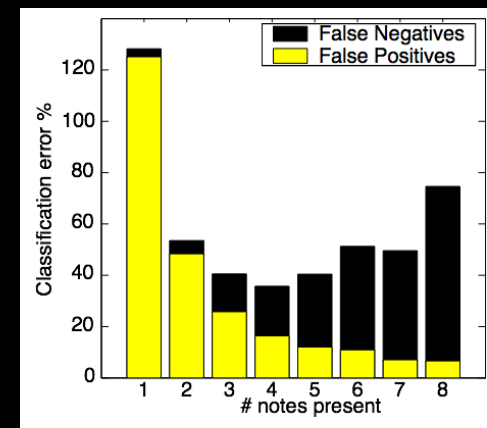IN THE CITY OF NEW YORK

# LabROSA

- Getting information from sound

# 1. Music Audio Analysis

- Trained classifiers for low-level information
  - notes, chords, beats, section boundaries
- E.g. Polyphonic transcription



- feature agnostic
- needs training data
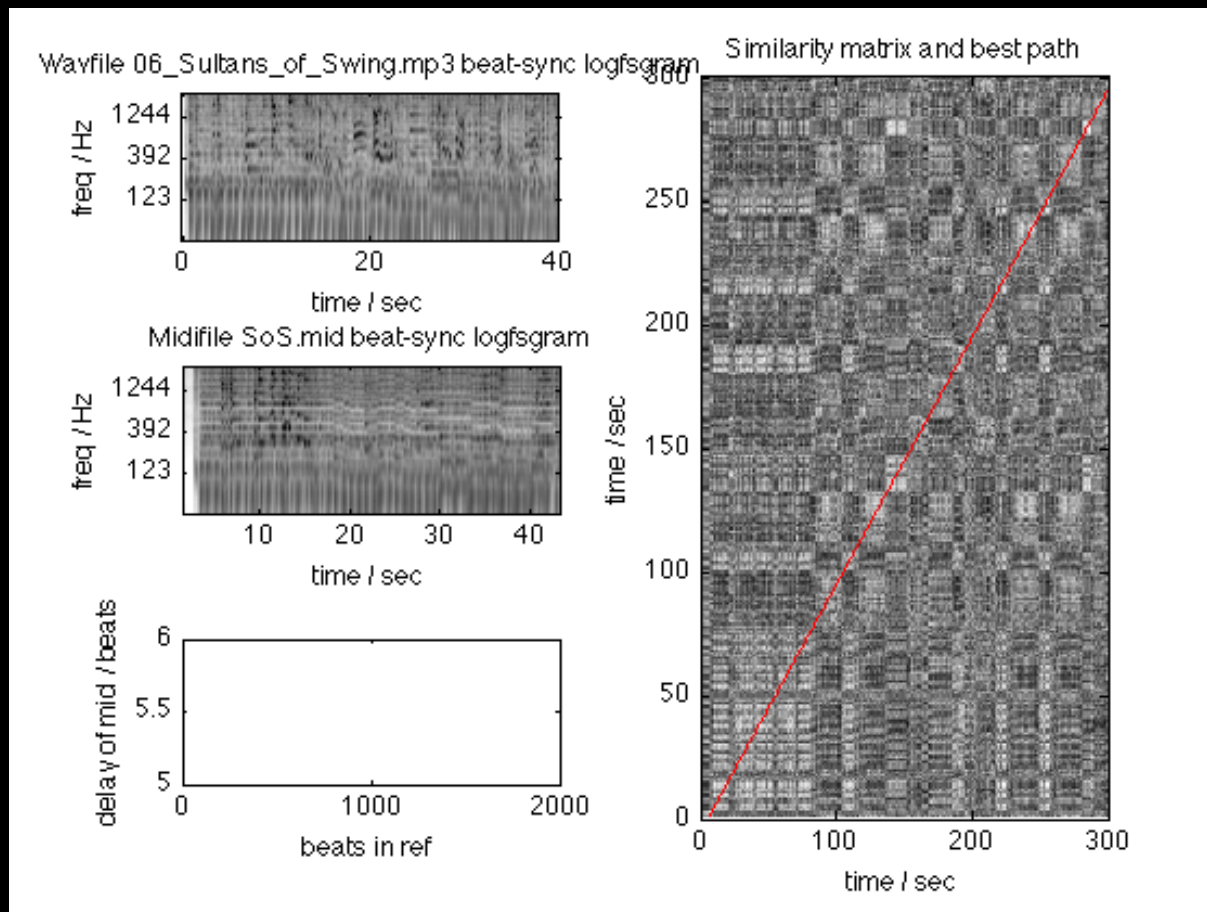
# Million Song Dataset

*Bertin-Mahieux*
*McFee*

- **Industrial-scale** database for music information research
- Many facets:
  - Echo Nest audio features + metadata
  - Echo Nest "taste profile" user-song-listen count
  - Second Hand Song covers
  - musiXmatch lyric BoW
  - last.fm tags
- **Now with audio?**
  - resolving artist / album / track / duration against what.cd

# MIDI-to-MSD

*Raffel*

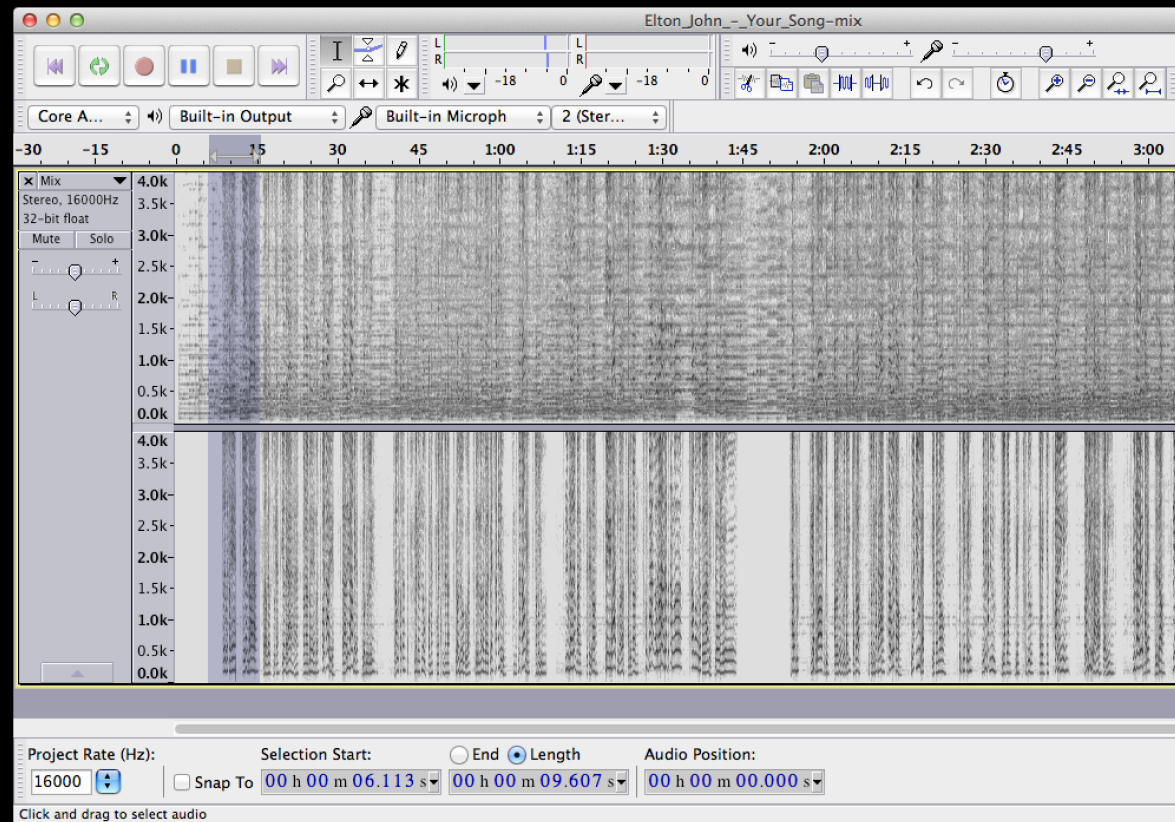- Aligned MIDI to Audio is a nice transcription



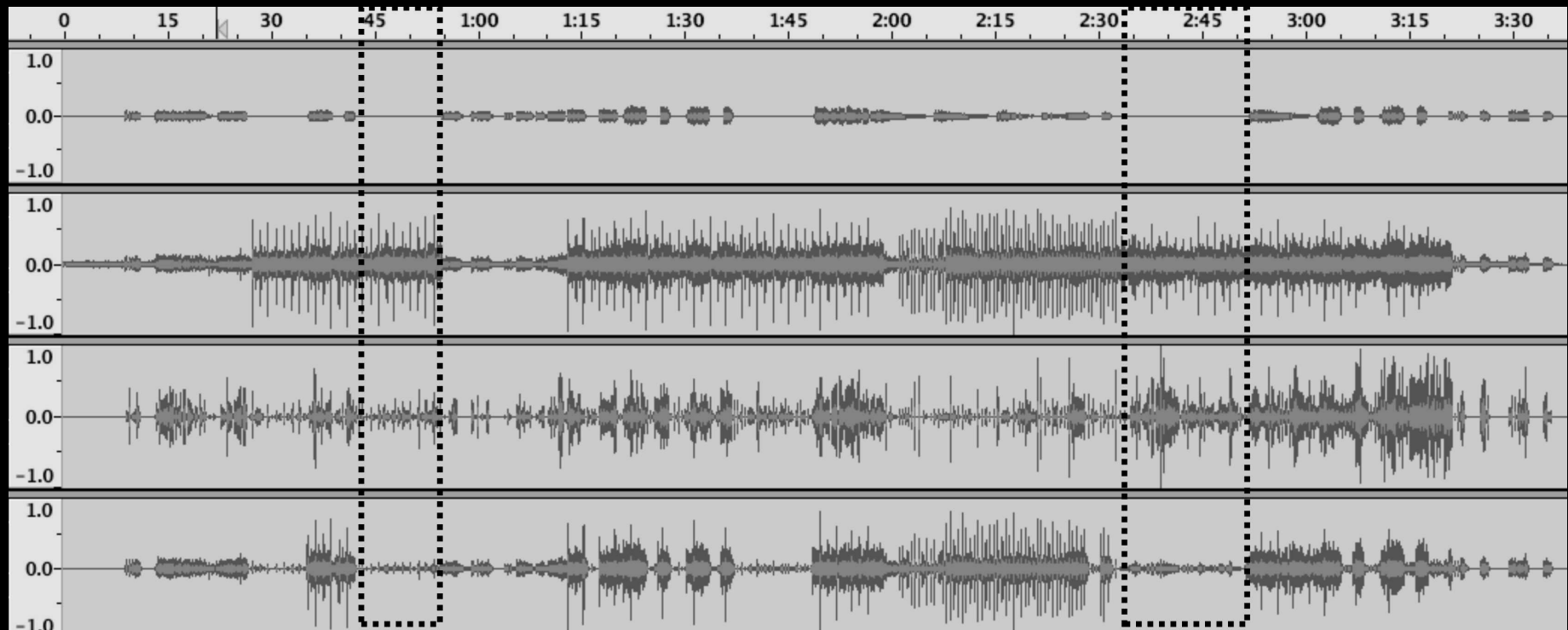- Can we find matches in large databases?

# Singing ASR

*McVicar*

- Speech recognition adapted to singing
  - needs aligned data
- Extensive work to line up scraped "acapellas" and full mix
  - including jumps!

# Block Structure RPCA

- **RPCA separates vocals and background based on low rank optimization**
  - single trade-off parameter
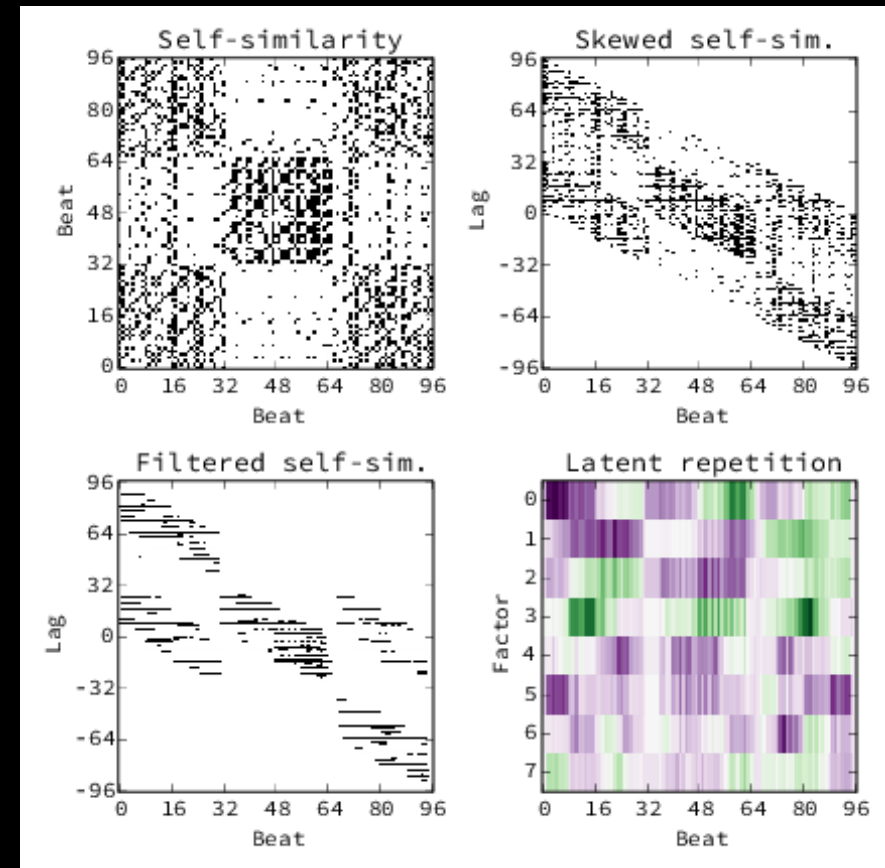  - adjust based on higher-level musical features?

# Ordinal LDA Segmentation

*McFee*

- Low-rank decomposition of skewed self-similarity to identify repeats
- Learned weighting of multiple factors to segment

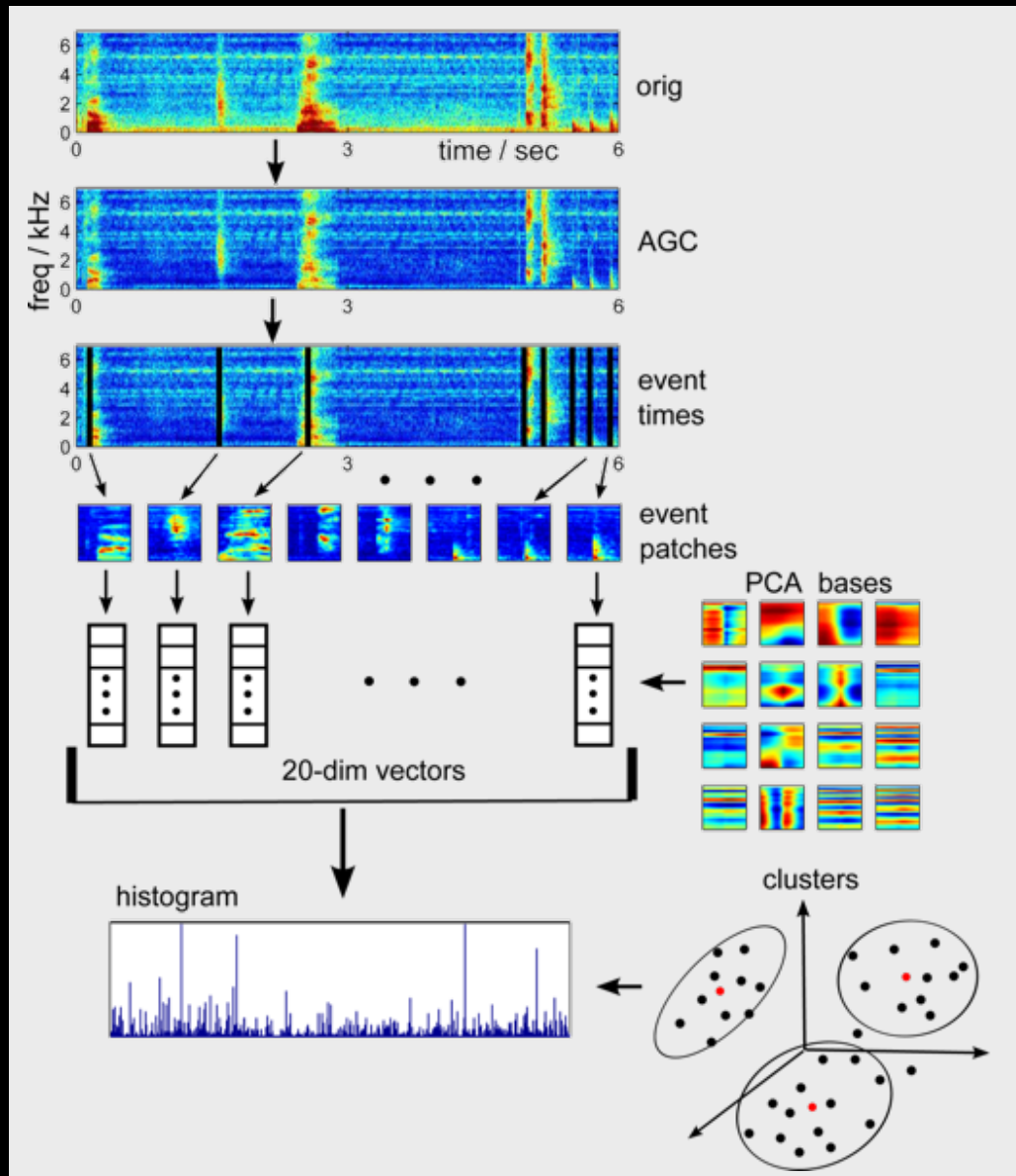  - Linear Discriminant Analysis between adjacent segments

# 2. Environmental Sound

- Extracting useful information from soundtracks

- e.g. TRECVID Multimedia Event Detection (MED)
  - "Making a Sandwich", "Getting a Vehicle Unstuck"
  - 100 examples, find matches in 100k videos
  - manual annotations for ~10 h



*E009 Getting a Vehicle Unstuck*

# Foreground Event Recognition

*Cotton, Ellis, Loui '11*



- Transients = foreground events?
- Onset detector finds energy bursts
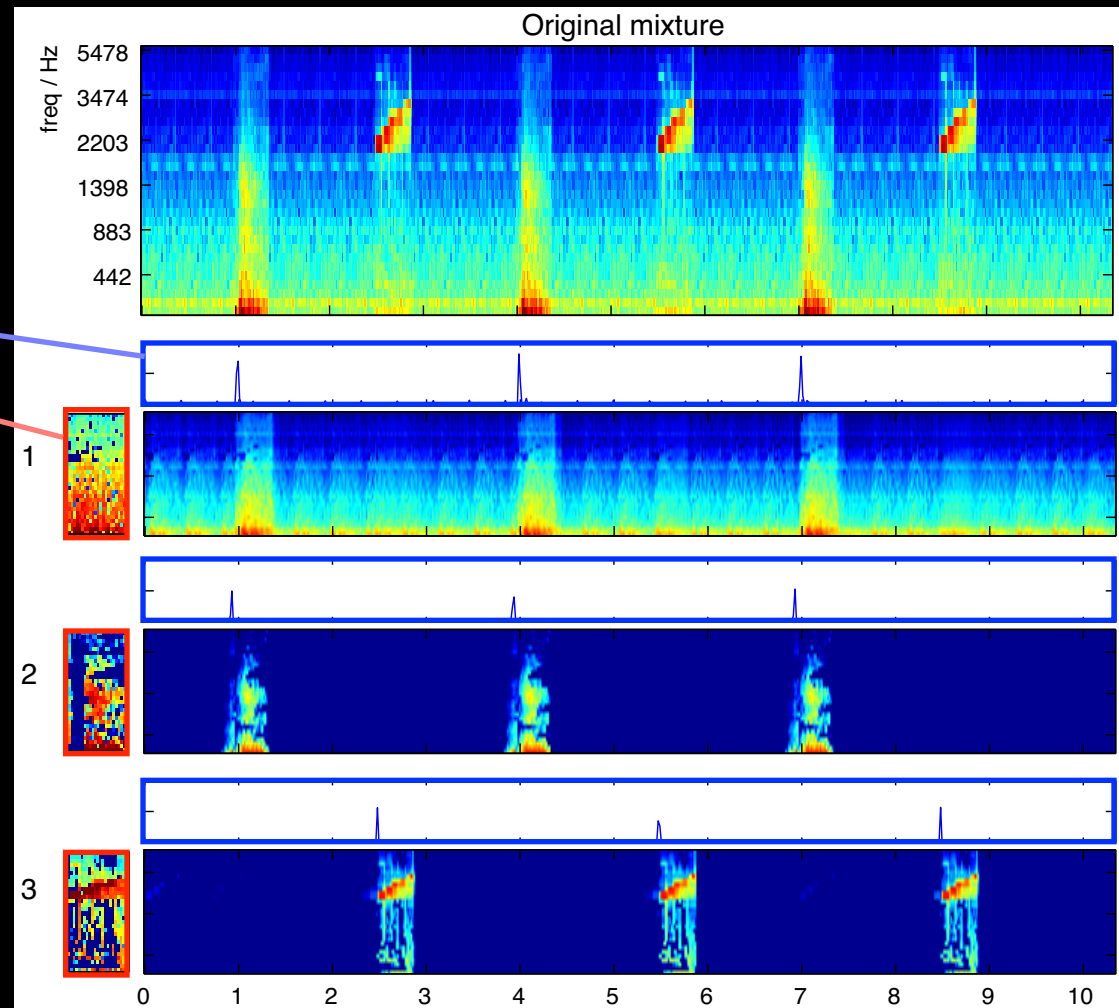  - best SNR
- PCA basis to represent each
  - 300 ms x auditory freq
- "bag of transients"

# NMF Transient Features

*Smaragdis & Brown '03*
*Abdallah & Plumbley '04*
*Virtanen '07*
*Cotton & Ellis' 11*

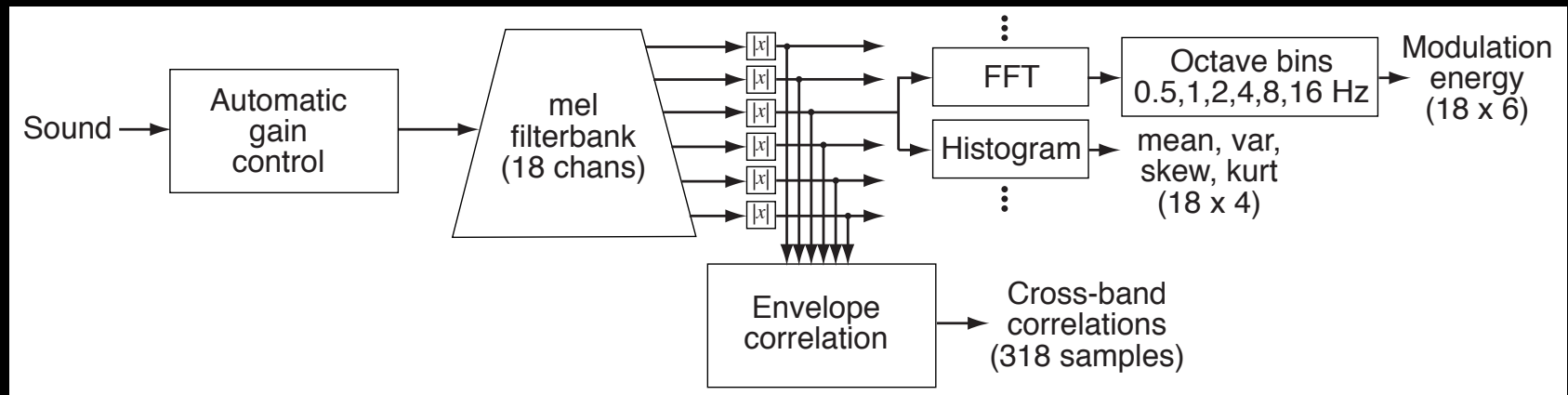- **Decompose spectrograms into templates + activation**

$$X = W \cdot H$$

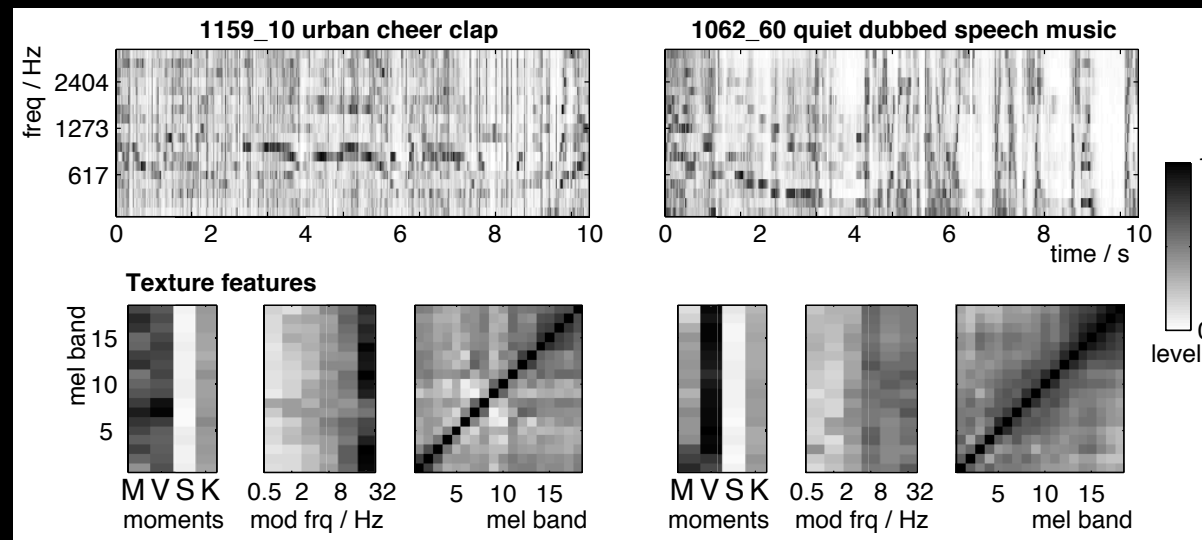- well-behaved gradient descent
- 2D patches
- sparsity control
- computation time…

# Background Retrieval

- Classify soundtracks by statistics of ambience
- E.g. Texture features



- Subband distributions
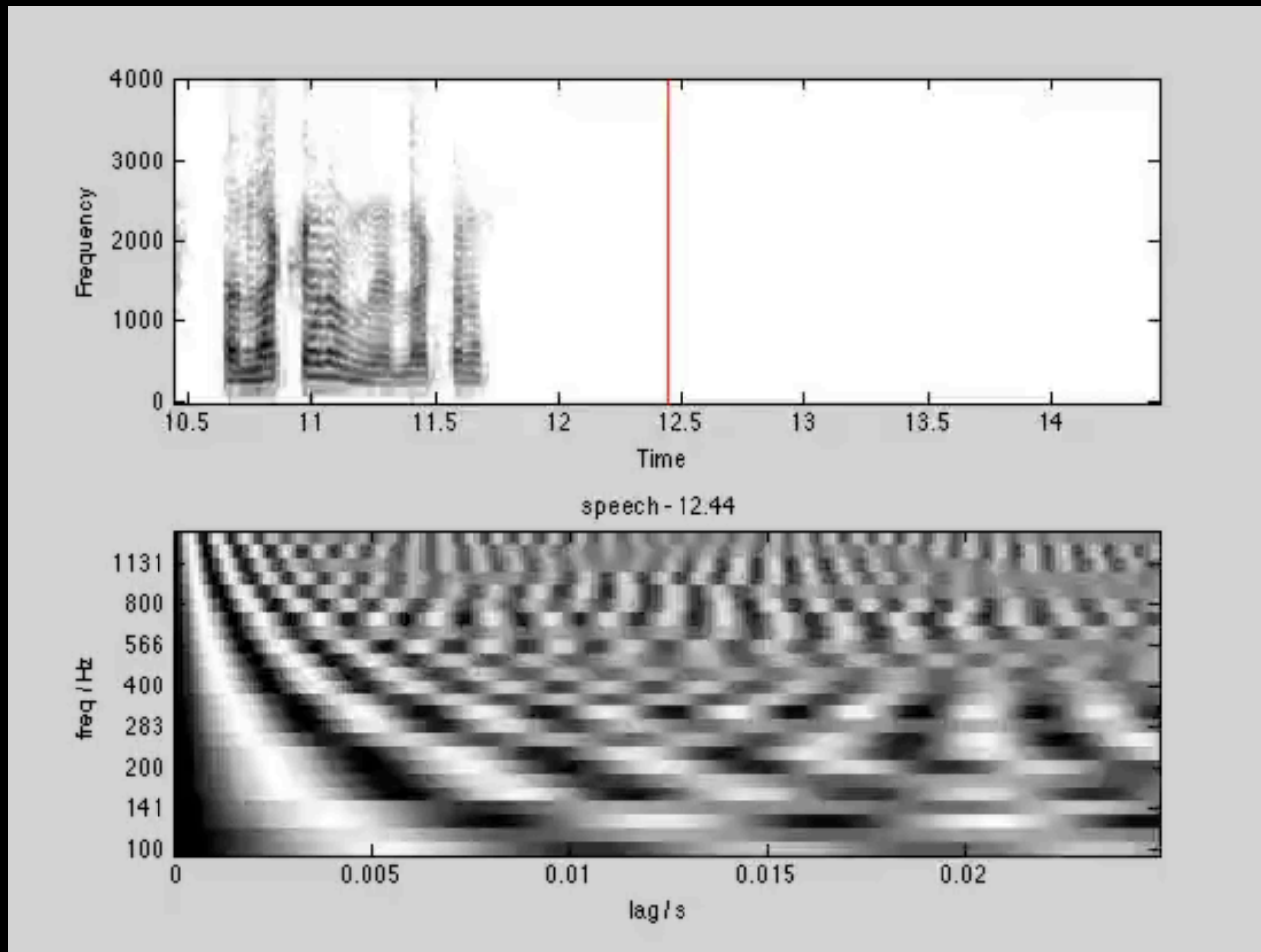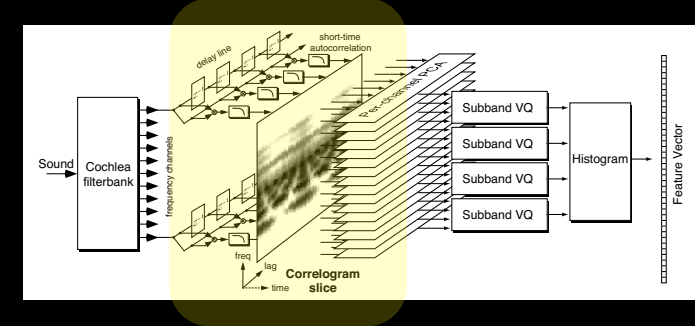- Envelope cross-corrs

# Auditory Model Features

- **Subband Autocorrelation PCA**
  - Simplified version of autocorrelogram
  - 10x faster than Lyon original
- **Capture fine time structure in multiple bands**
  - information lost in MFCCs

# Subband Autocorrelation



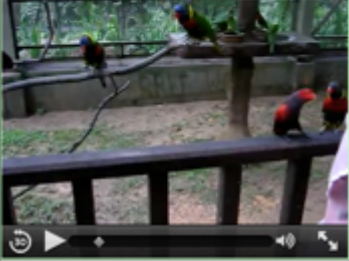- Autocorrelation **stabilizes** fine time structure
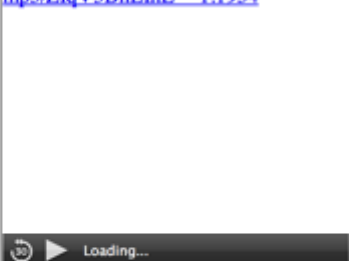


speech - 12.44

- 25 ms window, lags up to 25 ms
- calculated every 10 ms
- normalized to max (zero lag)

# Retrieval Examples
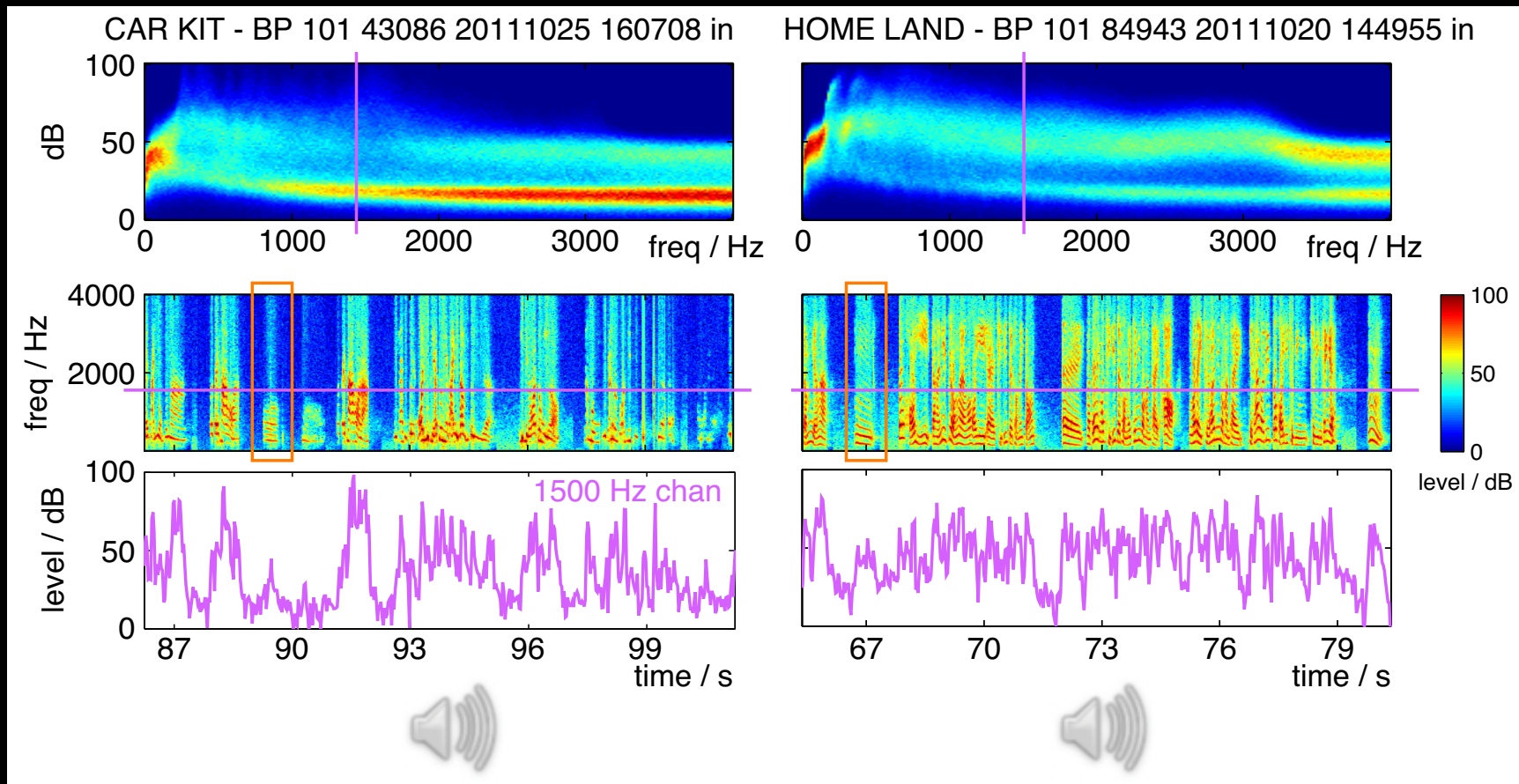
- High precision for in-domain top hits

# 3. Speech Enhancement

- Noisy speech scenarios
  - Ambient recording (background noise)
  - Communication channel (processing distortion)

# RPCA Enhancement
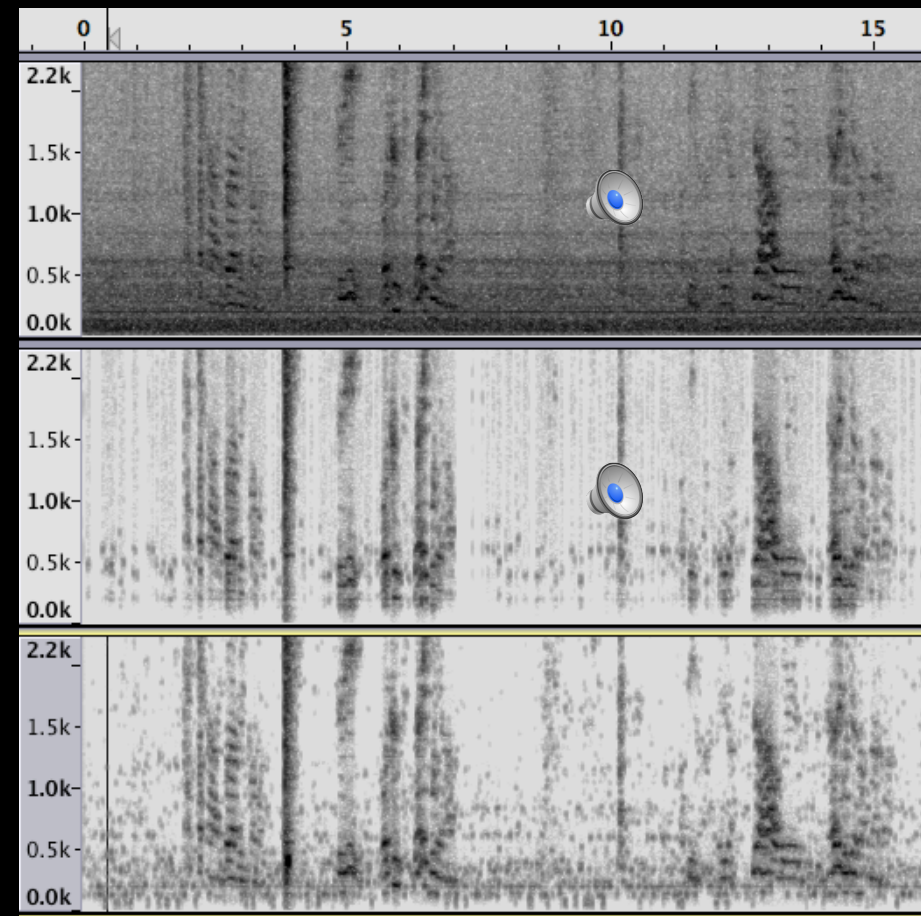
- Decompose spectrogram into sparse + low-rank

- Sparse activation H of dictionary W

$$\min_{H,L,S} \lambda_H \|H\|_1 + \lambda_L \|L\|_* + \lambda_S \|S\|_1$$

$$+ \mathcal{I}_+(H)$$

$$\text{s.t. } Y = WH + L + S$$

- ASR benefits:

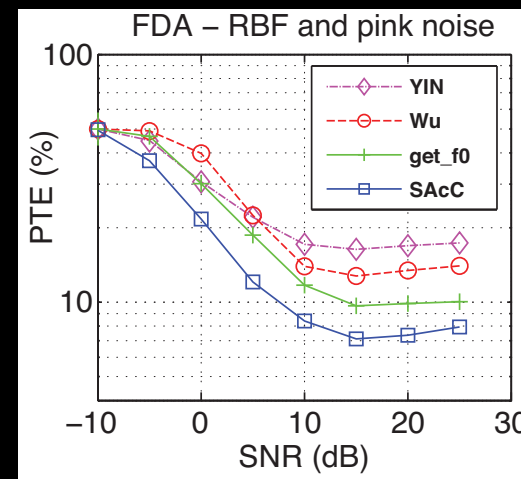| | C | S | D | I |
|---|---|---|---|---|
| Orig | 6.8 | 10.6 | 82.6 | 0.7 |
| RPCA | 10.8 | 36.5 | 52.7 | 0.5 |
| wie+RPCA | 10.4 | 40.1 | 49.5 | 2.1 |

# Classification Pitch Tracker

- SAcC: MLP trained on noisy speech with ground-truth pitch track targets



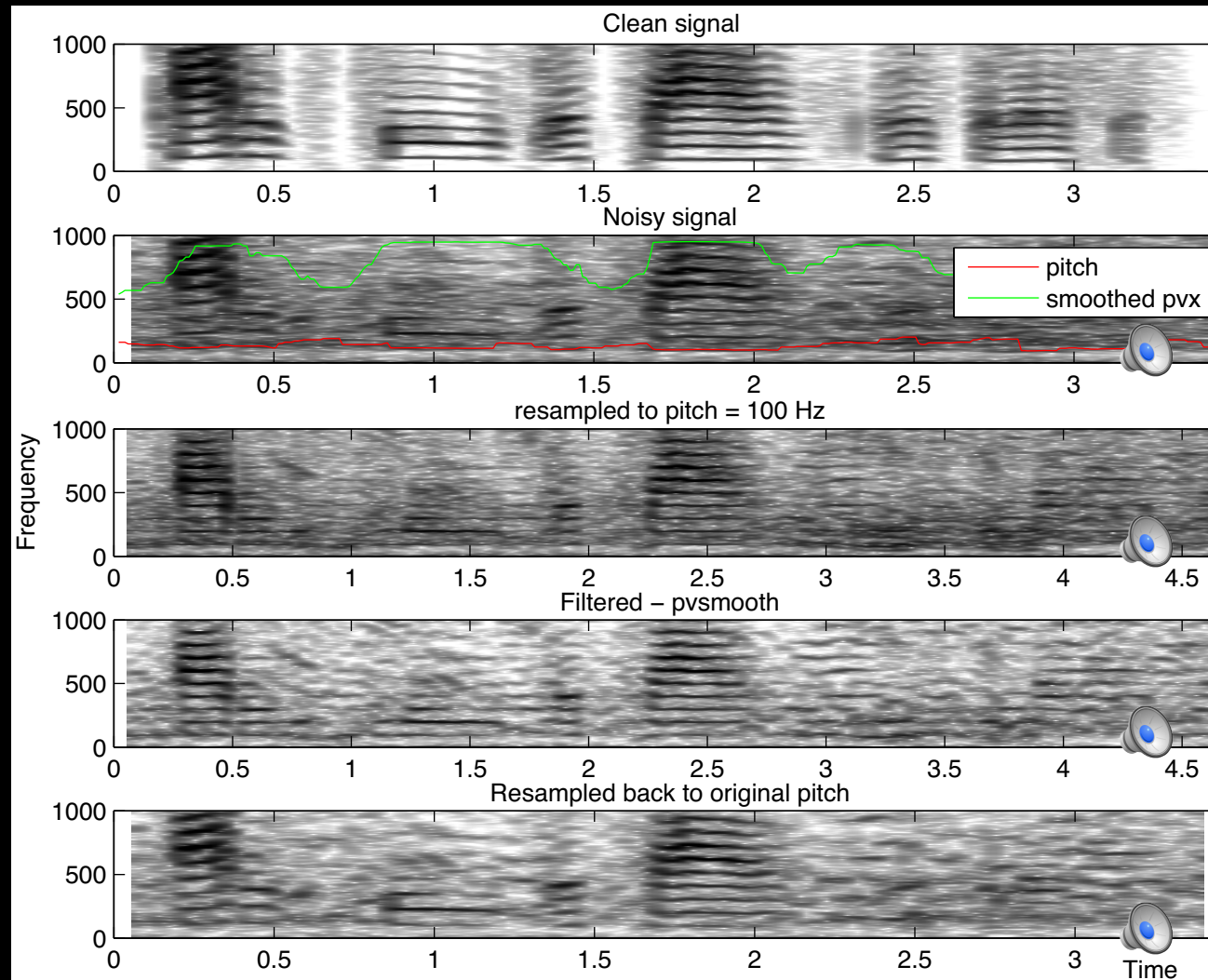- Large benefits for in-domain noisy speech

# Pitch-Normalized Enhancement

- ## Use noise-robust <span style="color:blue">pitch tracker</span> for <span style="color:red">enhancement</span>?

- Normalize voice pitch

- Fixed-pitch enhancement

- Reimpose pitch



Clean signal

Noisy signal
pitch
smoothed pvx

resampled to pitch = 100 Hz

Filtered – pvsmooth

Resampled back to original pitch

# Summary

- **Music**
  - transcription, segmentation, …
  - alignment for ground truth

- **Soundtracks**
  - foreground events, background ambience

- **Noisy Speech**
  - classification pitch tracking
  - spectrogram enhancement