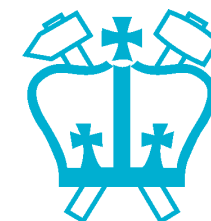# Extracting Information from Music Audio

Dan Ellis

Laboratory for Recognition and Organization of Speech and Audio
Dept. Electrical Engineering, Columbia University, NY USA

http://labrosa.ee.columbia.edu/

1. Learning Music
2. Melody Extraction
3. Music Similarity

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK
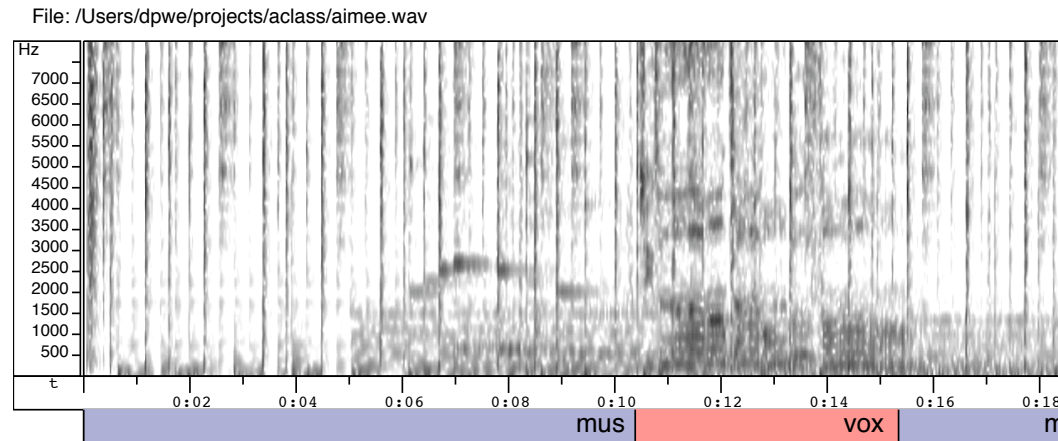
# I. Learning from Music

- **A lot of music data available**
  - e.g. 60G of MP3
    ≈ 1000 hr of audio, 15k tracks
- **What can we do with it?**
  - implicit definition of 'music'
- **Quality vs. quantity**
  - Speech recognition lesson:
    10x data, 1/10th annotation, twice as useful
- **Motivating Applications**
  - music similarity / classification
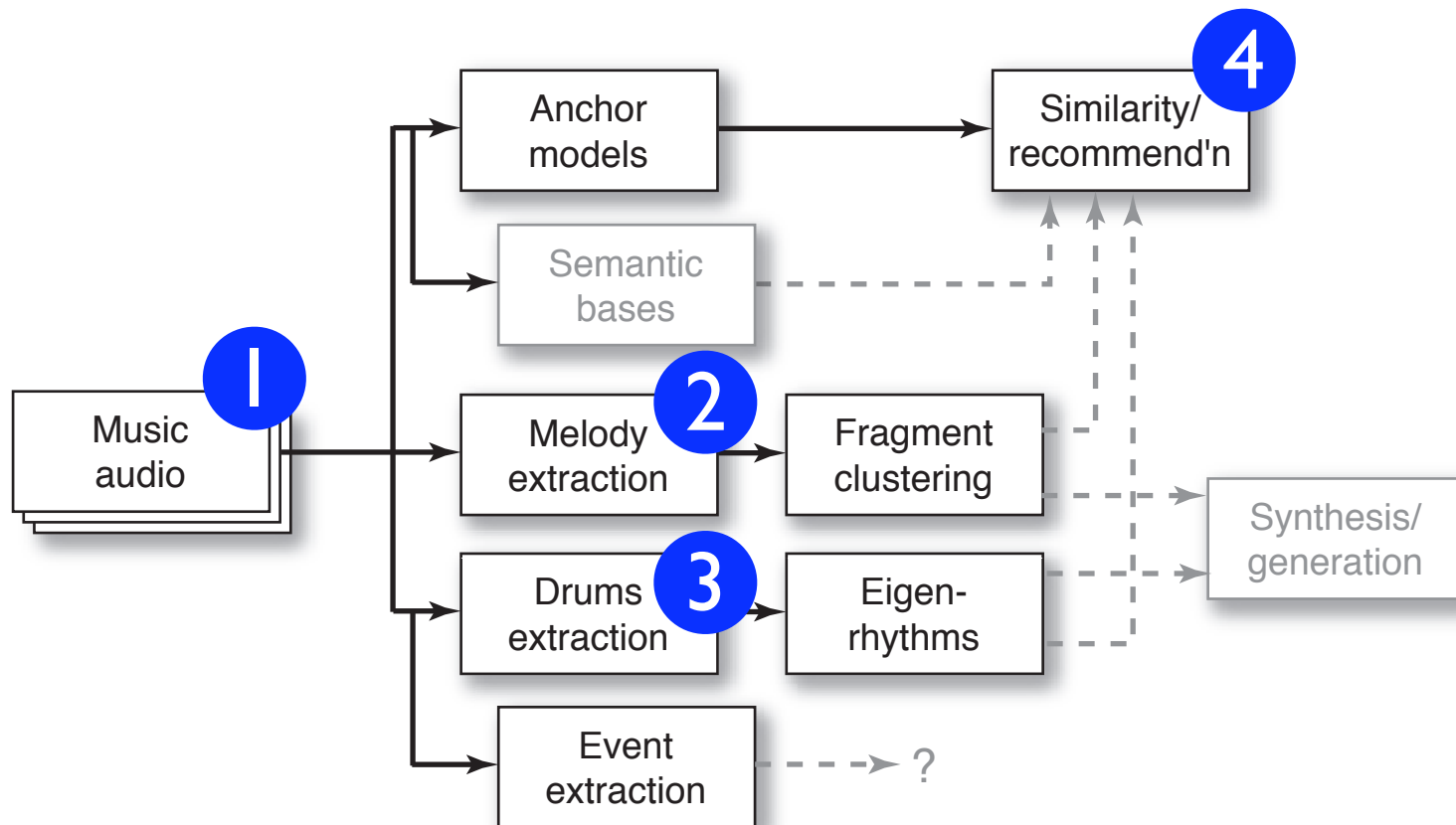  - computer (assisted) music generation
  - insight into music

1955 songs, 5.1 days, 7.90 GB

# Ground Truth Data

File: /Users/dpwe/projects/aclass/aimee.wav

- A lot of **unlabeled** music data available
  - ○ manual annotation is much rarer

- **Unsupervised structure discovery possible**
  - ○ .. but labels help to indicate what you want

- **Weak annotation sources**
  - ○ artist-level descriptions
  - ○ symbol sequences without timing (MIDI)
  - ○ errorful transcripts

- **Evaluation requires ground truth**
  - ○ limiting factor in Music IR evaluations?

LabROSA
Laboratory for the Recognition and
Organization of Speech and Audio
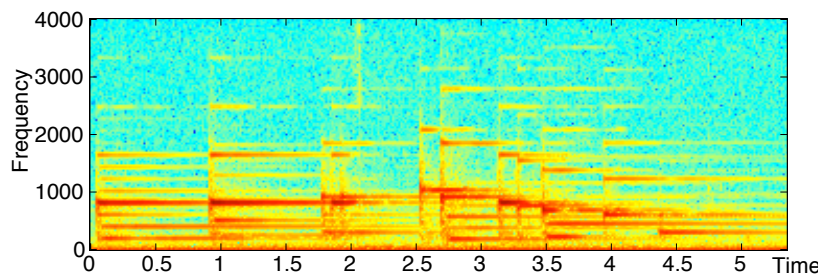
COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK
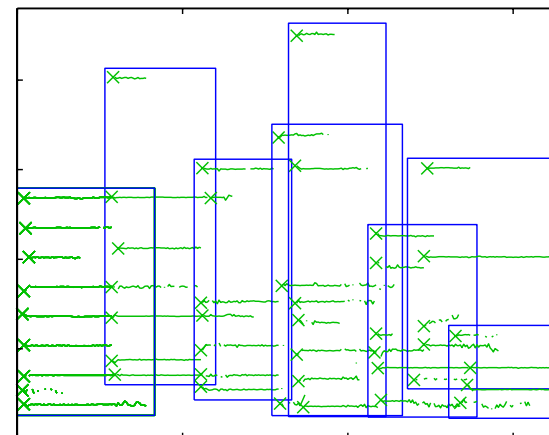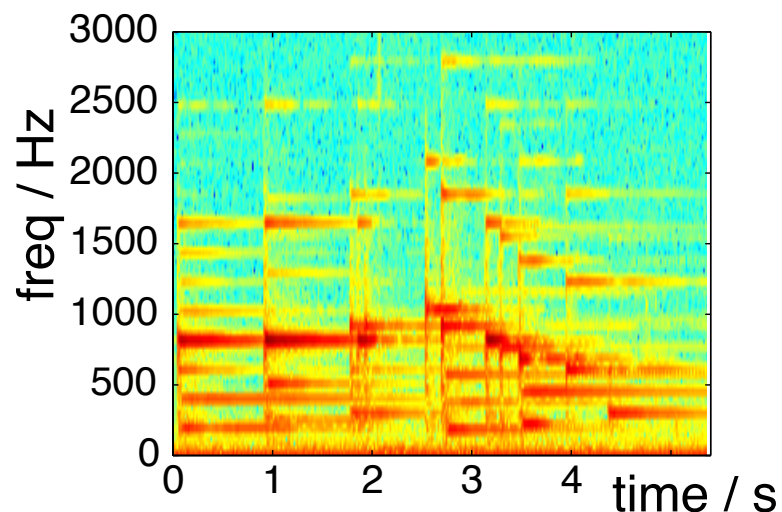
# Talk Roadmap

# 2. Melody Transcription

with Graham Poliner

- Audio → Score very desirable
  - for data compression, searching, learning
- **Full solution is elusive**
  - signal separation of overlapping voices
  - music constructed to frustrate!
- **Simplified problem:**
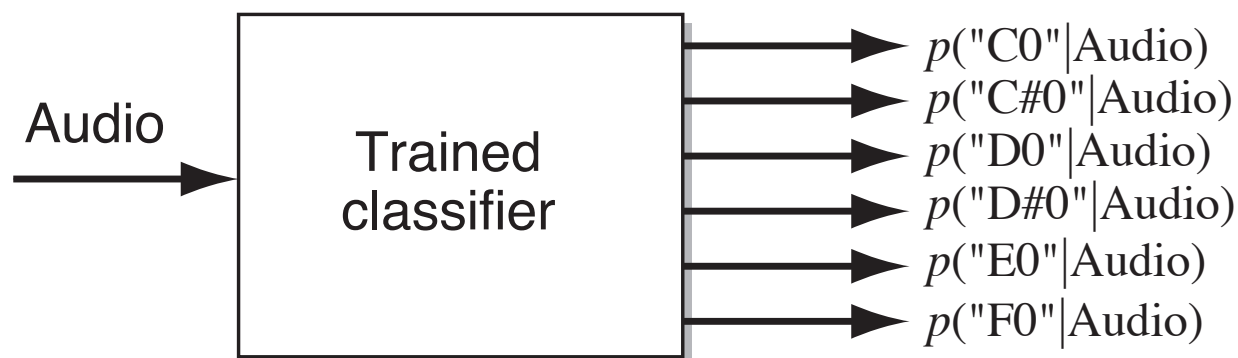  "Dominant Melody" at each time frame

# Conventional Transcription

- Pitched notes have harmonic spectra
  → transcribe by searching for harmonics
  - e.g. sinusoid modeling + grouping



- Explicit expert-derived knowledge

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

COLUMBIA UNIVERSITY
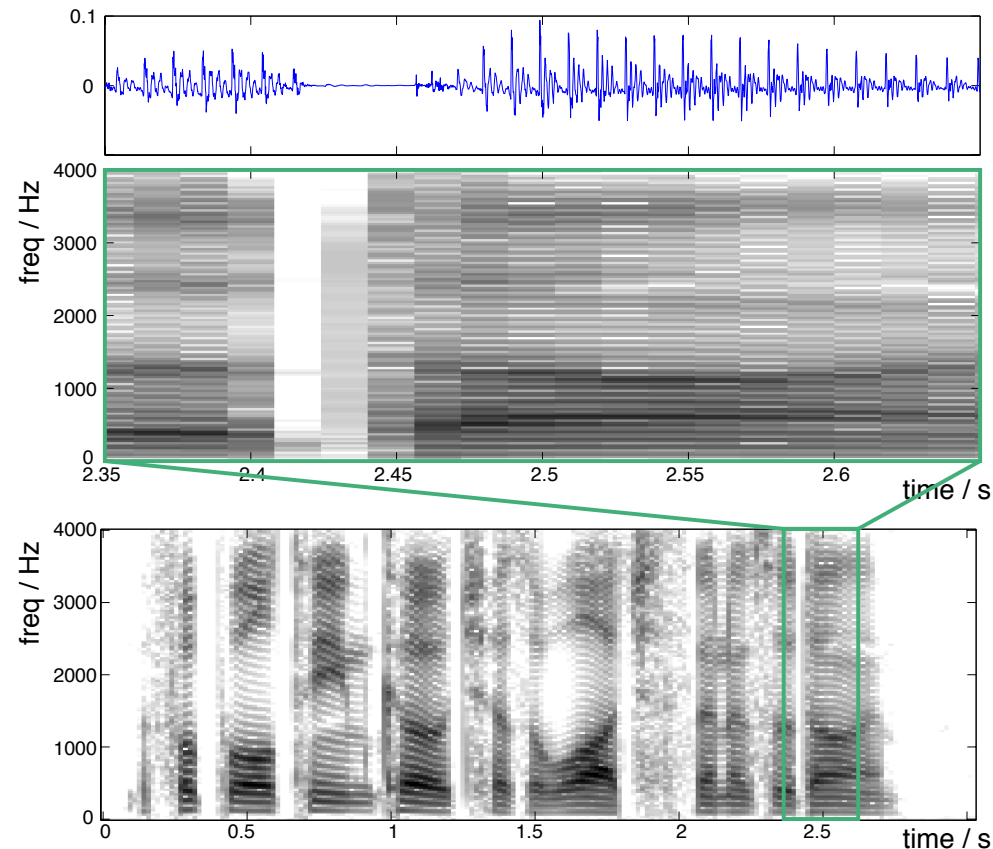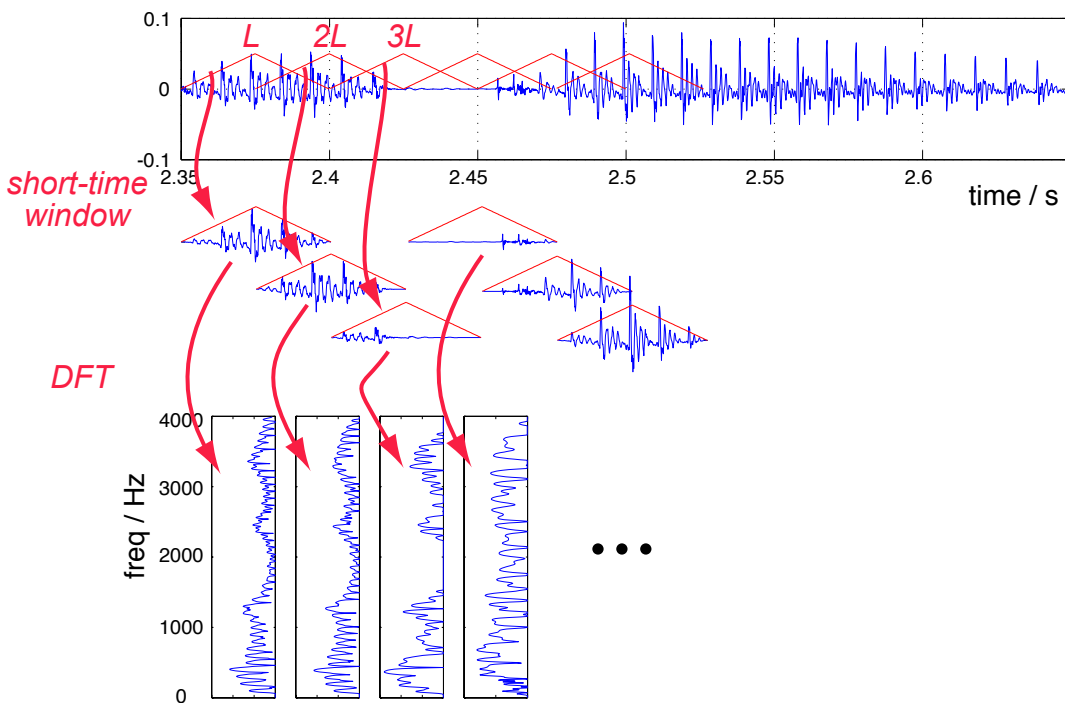IN THE CITY OF NEW YORK

# Transcription as Classification

- **Signal models** typically used for transcription
  - harmonic spectrum, superposition

- **But ... trade domain knowledge for data**
  - transcription as pure classification problem:

Audio → [ Trained classifier ] →
- $p(\text{"C0"}|\text{Audio})$
- $p(\text{"C\#0"}|\text{Audio})$
- $p(\text{"D0"}|\text{Audio})$
- $p(\text{"D\#0"}|\text{Audio})$
- $p(\text{"E0"}|\text{Audio})$
- $p(\text{"F0"}|\text{Audio})$

  - single N-way discrimination for "melody"
  - per-note classifiers for polyphonic transcription

Lab ROSA

Laboratory for the Recognition and
Organization of Speech and Audio

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Melody Transcription Features

- Short-time Fourier Transform Magnitude (Spectrogram)
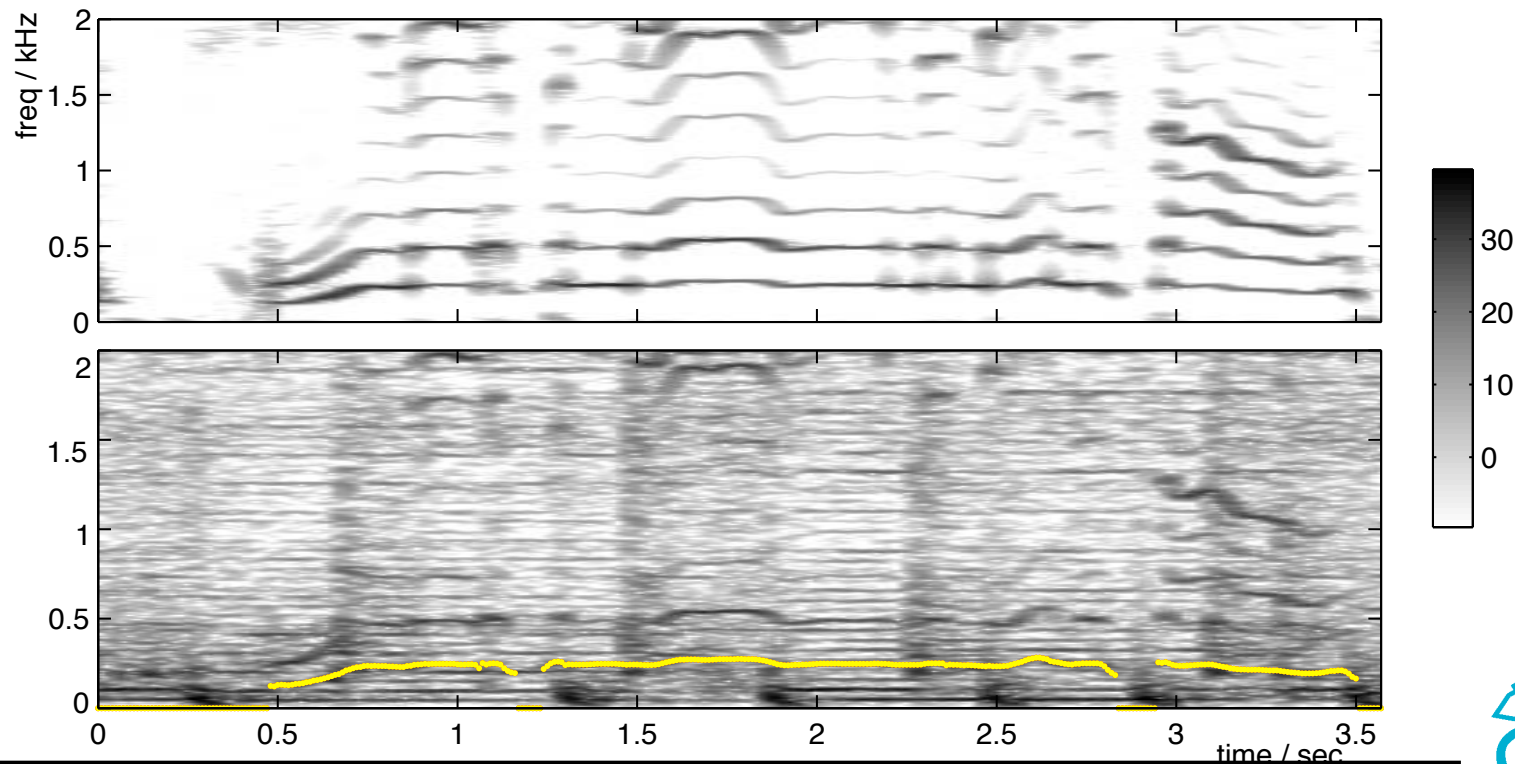


- Standardize over 50 pt frequency window

# Training Data

- Need {data, label} pairs for classifier training
- Sources:
  - pre-mixing multitrack recordings + hand-labeling?
  - synthetic music (MIDI)  + forced-alignment?
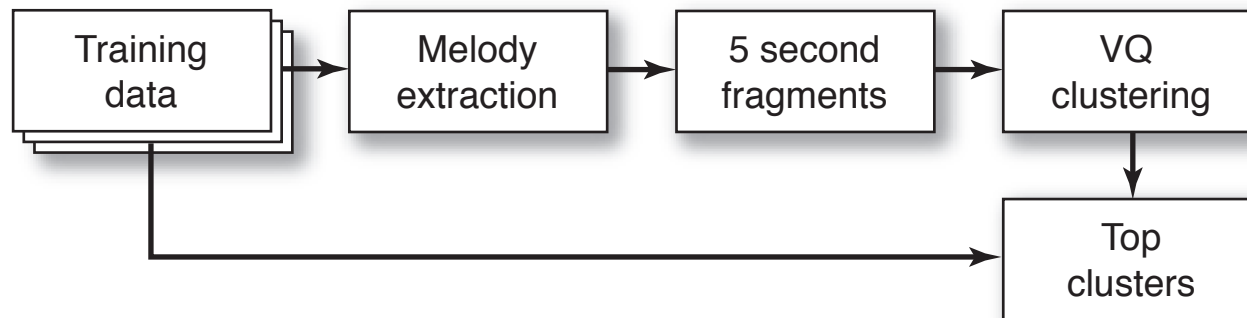
# Melody Transcription Results

- ## Trained on 17 examples
  - ○ .. plus transpositions out to +/- 6 semitones
  - ○ SMO SVM (Weka)
- ## Tested on ISMIR MIREX 2005 set
  - ○ includes foreground/background detection

| Rank | Participant | Overall Accuracy | Voicing $d'$ | Raw Pitch | Raw Chroma | Runtime / s |
|------|-------------|------------------|--------------|-----------|------------|-------------|
| 1 | Dressler | **71.4%** | **1.85** | 68.1% | 71.4% | 32 |
| 2 | Ryynänen | 64.3% | 1.56 | **68.6%** | **74.1%** | 10970 |
| 3 | Paiva 2 | 61.1% | 1.22 | 58.5% | 62.0% | 45618 |
| 3 | Poliner | 61.1% | 1.56 | 67.3% | 73.4% | 5471 |
| 5 | Marolt | 59.5% | 1.06 | 60.1% | 67.1% | 12461 |
| 6 | Paiva 1 | 57.8% | 0.83 | 62.7% | 66.7% | 44312 |
| 7 | Goto | 49.9%* | 0.59* | 65.8% | 71.8% | 211 |
| 8 | Vincent 1 | 47.9%* | 0.23* | 59.8% | 67.6% | ? |
| 9 | Vincent 2 | 46.4%* | 0.86* | 59.6% | 71.1% | 251 |
| 10 | Brossier | 3.2%* † | 0.14 * † | 3.9% † | 8.1% † | 41 |

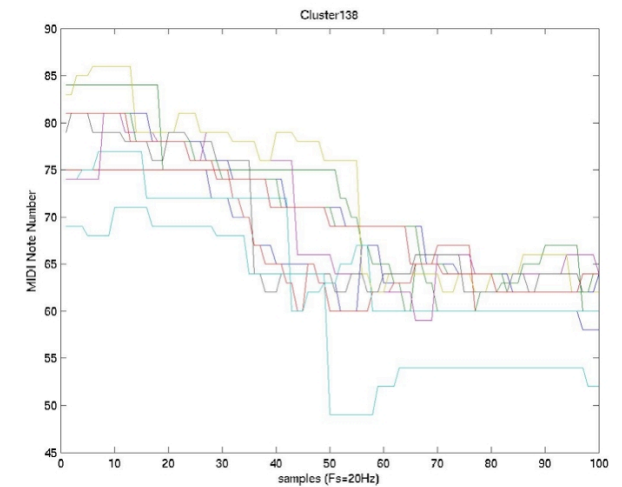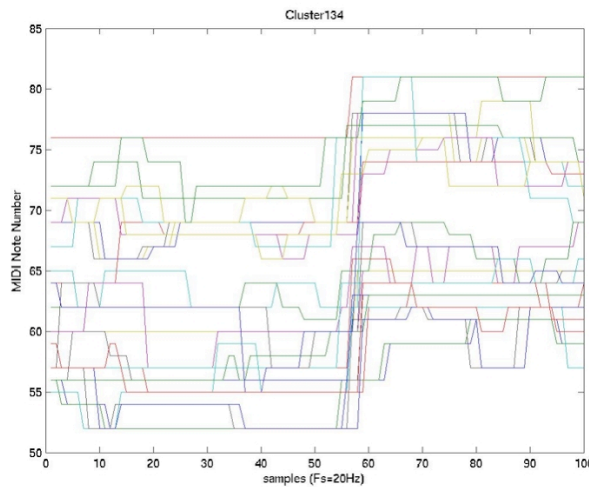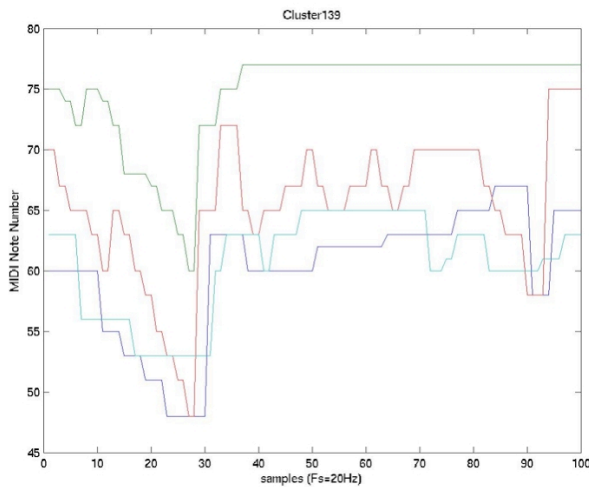- ○ Example...

# Melody Clustering

- **Goal: Find 'fragments' that recur in melodies**
  - .. across large music database
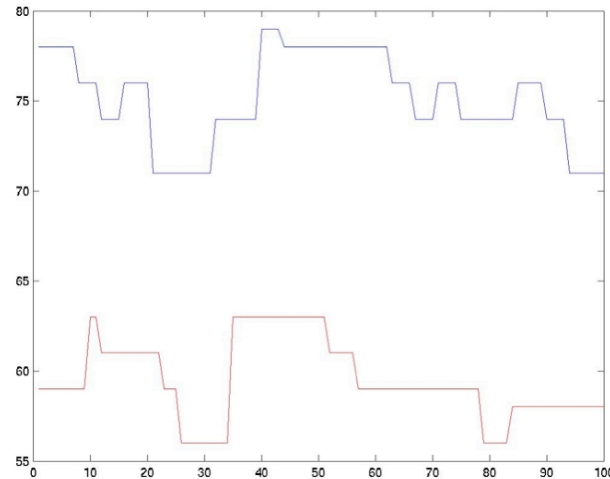  - .. trade data for model sophistication



- **Data sources**
  - pitch tracker, or MIDI training data
- **Melody fragment representation**
  - DCT(1:20) - removes average, smoothes detail

# Melody clustering results
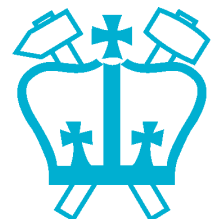
- Clusters match underlying contour:



- Some interesting matches:
  - e.g. Pink + Nsync

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# 3. Music Similarity

with Mike Mandel and Adam Berenzweig

- **Can we predict which songs "sound alike" to a listener?**
  - .. based on the audio waveforms?
  - many aspects to subjective similarity
- **Applications**
  - query-by-example
  - automatic playlist generation
  - discovering new music
- **Problems**
  - the right representation
  - modeling individual similarity

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Music Similarity Features

- Need "timbral" features:
Mel-Frequency Cepstral Coeffs (MFCCs)

  ○ auditory-like frequency warping

  ○ log-domain

  ○ discrete cosine transform orthogonalization

**Spectrogram**
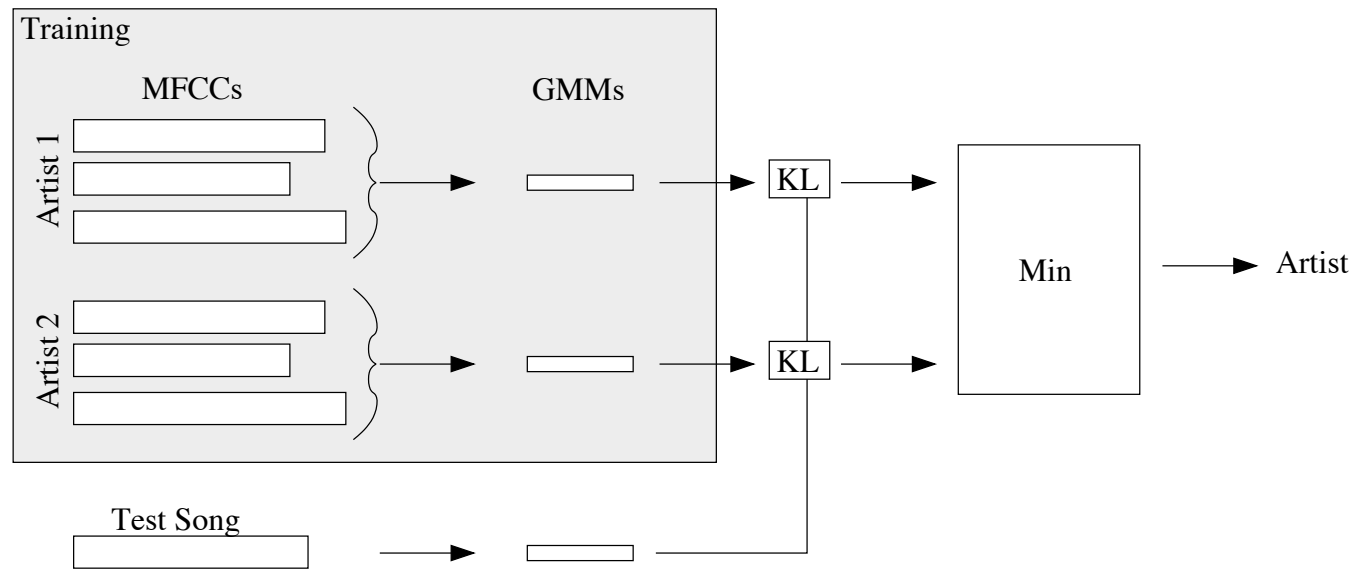
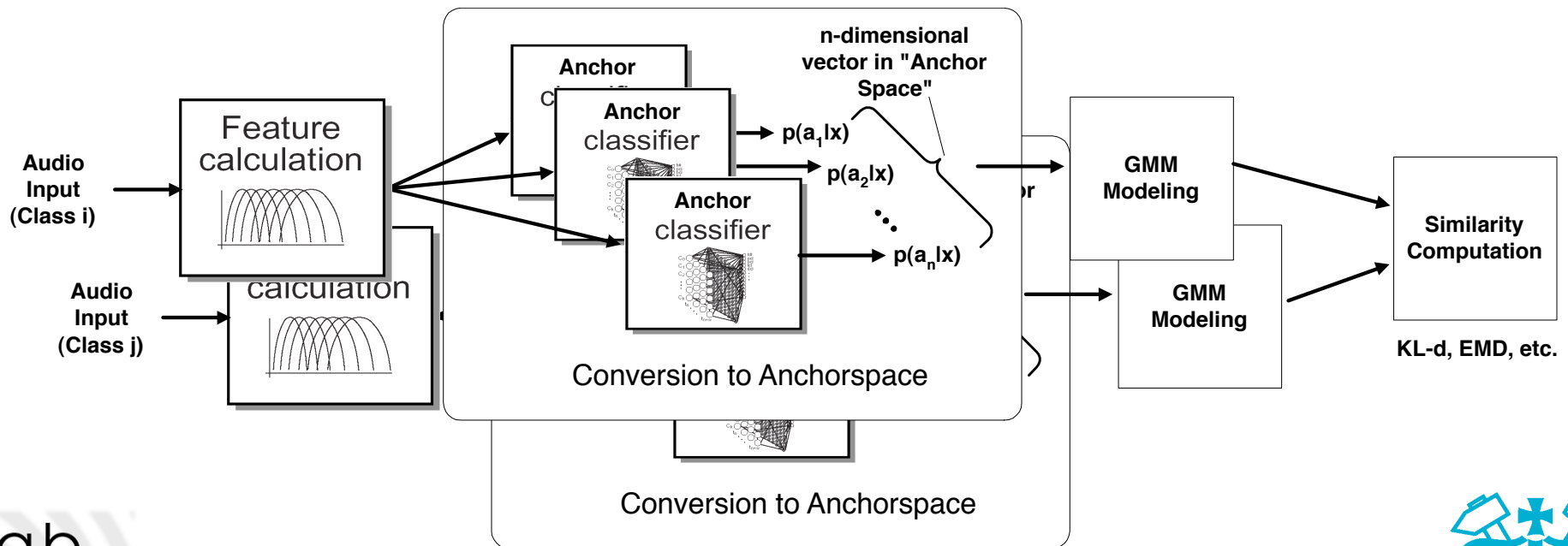**Mel-frequency Spectrogram**

**Mel-Frequency Cepstral Coefficients**

LabROSA
Laboratory for the Recognition and
Organization of Speech and Audio

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Timbral Music Similarity

- **Measure similarity of feature distribution**
  - i.e. collapse across time to get density $p(x_i)$
  - compare by e.g. KL divergence

- **e.g. Artist Identification**
  - learn artist model $p(x_i | \text{artist } X)$ (e.g. as GMM)
  - classify unknown song to closest model

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# "Anchor Space"

- ## Acoustic features describe each song
  - o .. but from a signal, not a perceptual, perspective
  - o .. and not the differences between songs
- ## Use genre classifiers to define new space
  - o prototype genres are "anchors"

# Anchor Space

- Frame-by-frame high-level categorizations
  - compare to raw features?



  - properties in distributions? dynamics?

# 'Playola' Similarity Browser

# Ground-truth data

- **Hard to evaluate Playola's 'accuracy'**
  - ○ user tests...
  - ○ ground truth?

- **"Musicseer" online survey:**
  - ○ ran for 9 months in 2002
  - ○ > 1,000 users, > 20k judgments
  - ○ http://labrosa.ee.columbia.edu/projects/musicsim/

Which artist is most similar to:
**Janet Jackson?**

1. R. Kelly
2. Paula Abdul
3. Aaliyah
4. Milli Vanilli
5. En Vogue
6. Kansas
7. Garbage
8. Pink
9. Christina Aguilera

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Evaluation

- ## Compare Classifier measures against Musicseer subjective results
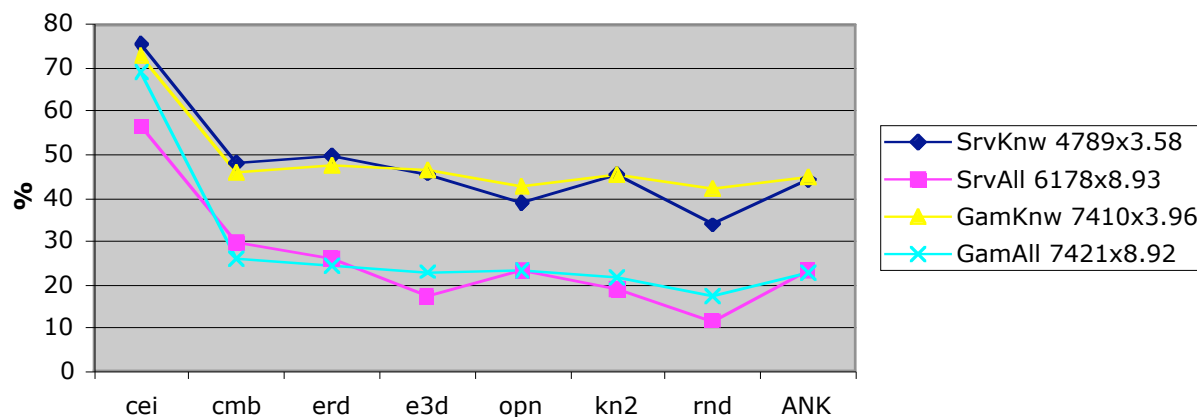
  - ○ "triplet" agreement percentage
  - ○ Top-N ranking agreement score:

$$s_i = \sum_{r=1}^{N} \alpha_r^r \alpha_c^{k_r} \qquad \alpha_r = \left(\frac{1}{2}\right)^{\frac{1}{3}} \qquad \alpha_c = \alpha_r^2$$

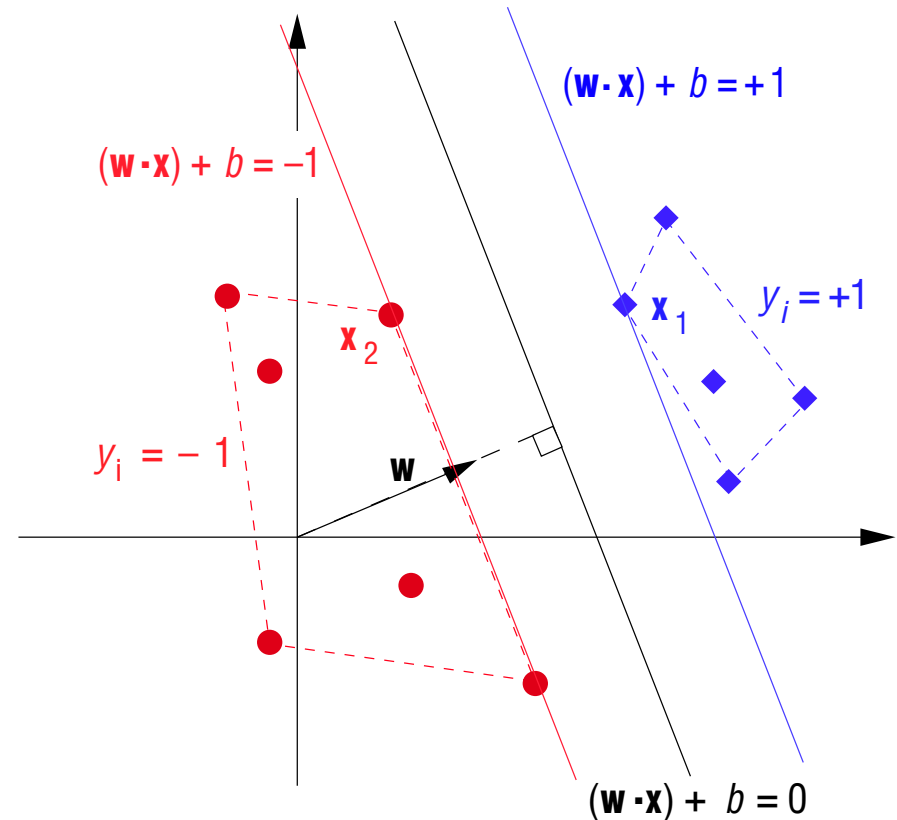  - ○ First-place agreement percentage
    - simple significance test

# Using SVMs for Artist ID

- Support Vector Machines (SVMs) find hyperplanes in a high-dimensional space
  - relies only on matrix of distances between points
  - much 'smarter' than nearest-neighbor/overlap
  - want diversity of reference vectors...

$(\mathbf{w} \cdot \mathbf{x}) + b = +1$

$(\mathbf{w} \cdot \mathbf{x}) + b = -1$

$y_i = +1$

$\mathbf{x}_1$

$\mathbf{x}_2$

$y_i = -1$

$\mathbf{w}$

$(\mathbf{w} \cdot \mathbf{x}) + b = 0$

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

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Song-Level SVM Artist ID

- Instead of one model per artist/genre, use *every* training song as an 'anchor'
  - then SVM finds best support for each artist

# Artist ID Results

- ISMIR/MIREX 2005 also evaluated Artist ID
- 148 artists, 1800 files (split train/test) from 'uspop2002'
- Song-level SVM clearly dominates
  - using only MFCCs!

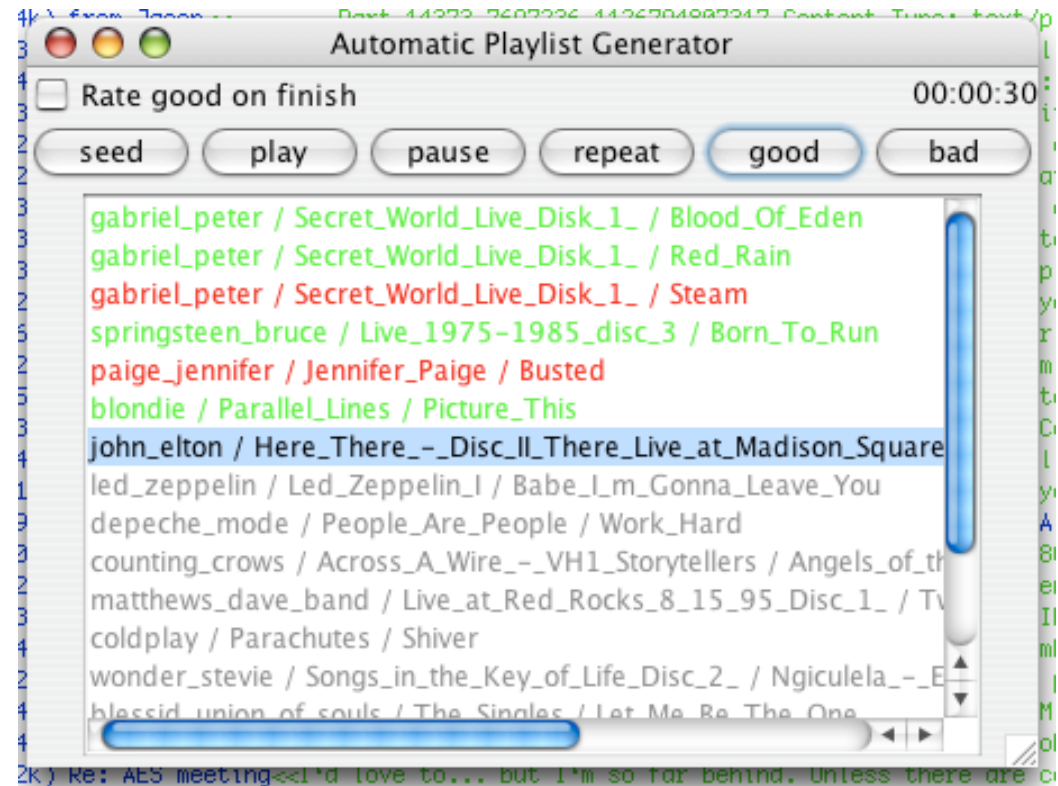MIREX 05 Audio Artist (USPOP2002)

| Rank | Participant | Raw Accuracy | Normalized | Runtime / s |
|------|-------------|--------------|------------|-------------|
| 1 | Mandel | **68.3%** | **68.0%** | 10240 |
| 2 | Bergstra | 59.9% | 60.9% | 86400 |
| 3 | Pampalk | 56.2% | 56.0% | 4321 |
| 4 | West | 41.0% | 41.0% | 26871 |
| 5 | Tzanetakis | 28.6% | 28.5% | 2443 |
| 6 | Logan | 14.8% | 14.8% | ? |
| 7 | Lidy | Did not complete | | |

LabROSA
Laboratory for the Recognition and Organization of Speech and Audio
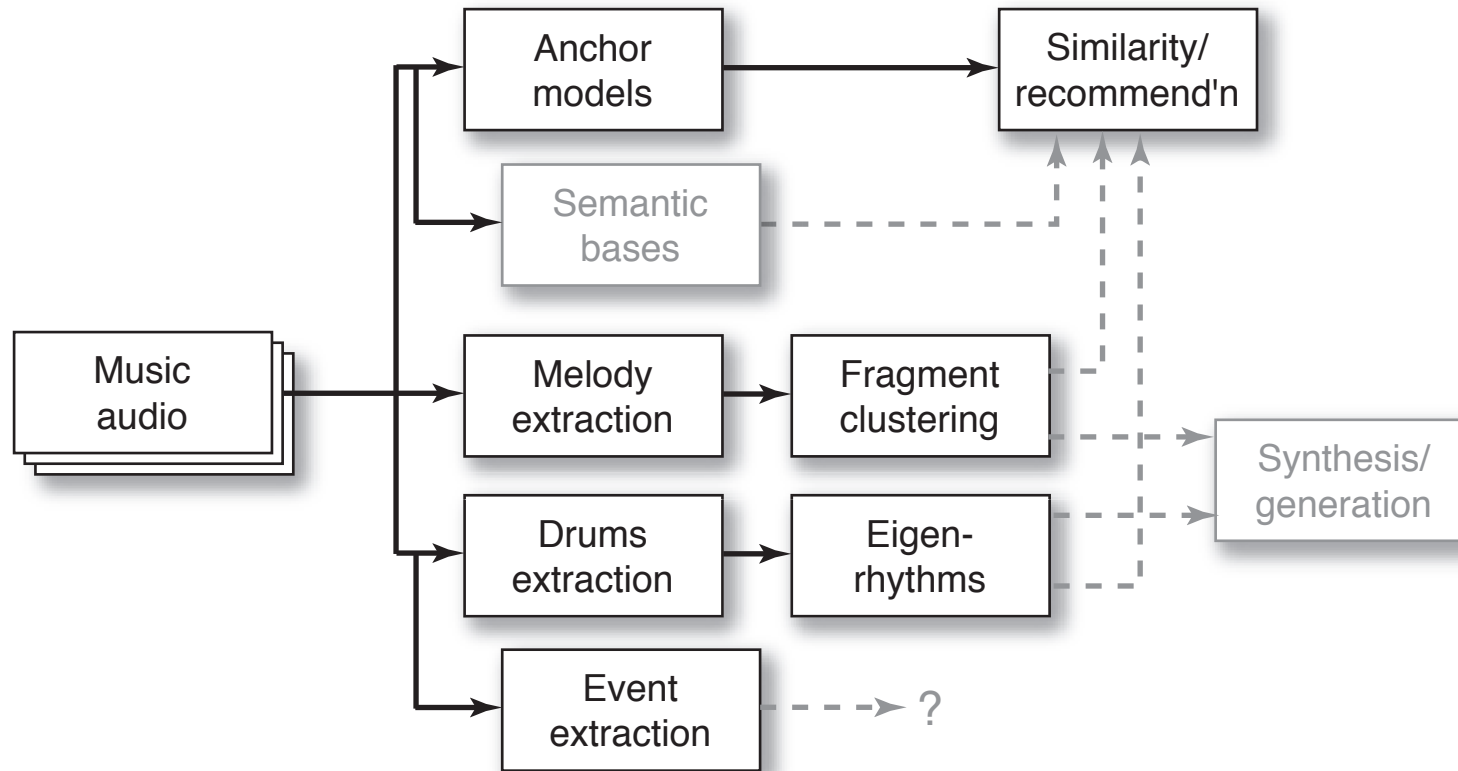
COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# Playlist Generation

- SVMs are well suited to "active learning"
  - solicit labels on items closest to current boundary

- Automatic player with "skip" = Ground truth data collection
  - active-SVM automatic playlist generation

# Conclusions



- Lots of data
  + noisy transcription
  + weak clustering
  ⇒ musical insights?