

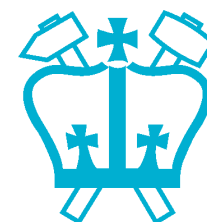
Audio & Music Research at LabROSA

Dan Ellis

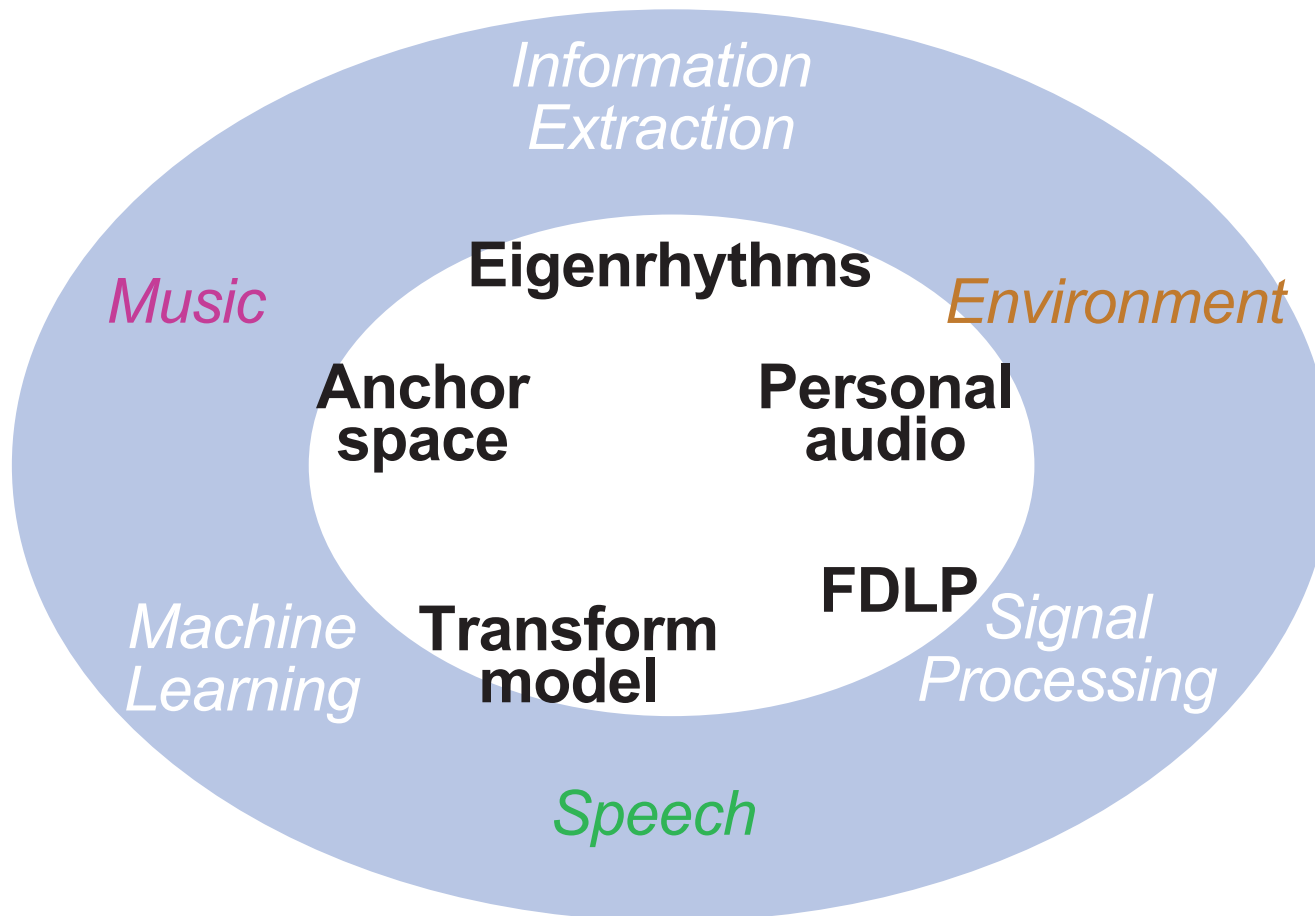
Laboratory for Recognition and Organization of Speech and Audio
Dept. Electrical Eng., Columbia Univ., NY USA

dpwe@ee.columbia.edu <http://labrosa.ee.columbia.edu/>

1. **Eigenrhythms**: representing drum tracks
2. **Frequency-Domain Linear Prediction**
3. **Anchor-Space** Music Similarity Browsing
4. **Transformation**-based generative models
5. Analyzing **'personal audio'** recordings



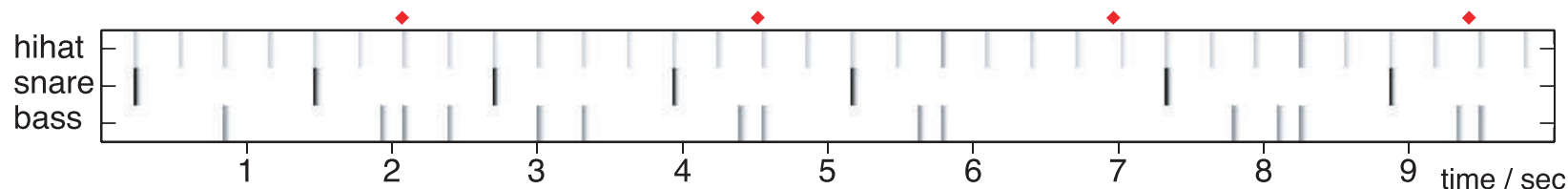
LabROSA Projects Overview



I. Eigenrhythms: Drum Pattern Space

with John Arroyo

- Pop songs built on repeating “drum loop”
 - bass drum, snare, hi-hat
 - small variations on a few basic patterns



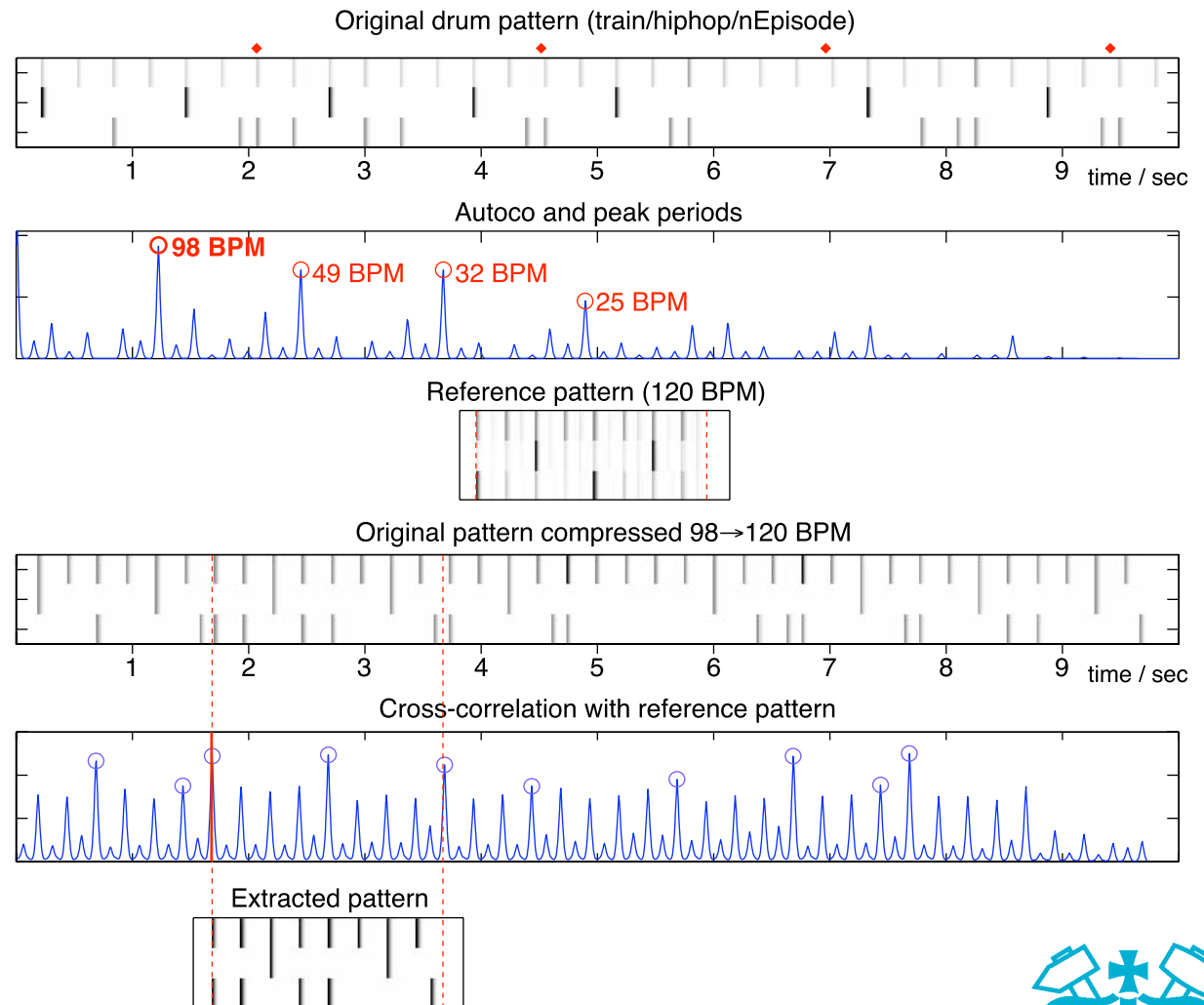
-
- **Eigen-analysis (PCA)** to capture variations?
 - by analyzing lots of (MIDI) data
- **Applications**
 - music categorization
 - “beat box” synthesis

Aligning the Data

- Need to align patterns prior to PCA...

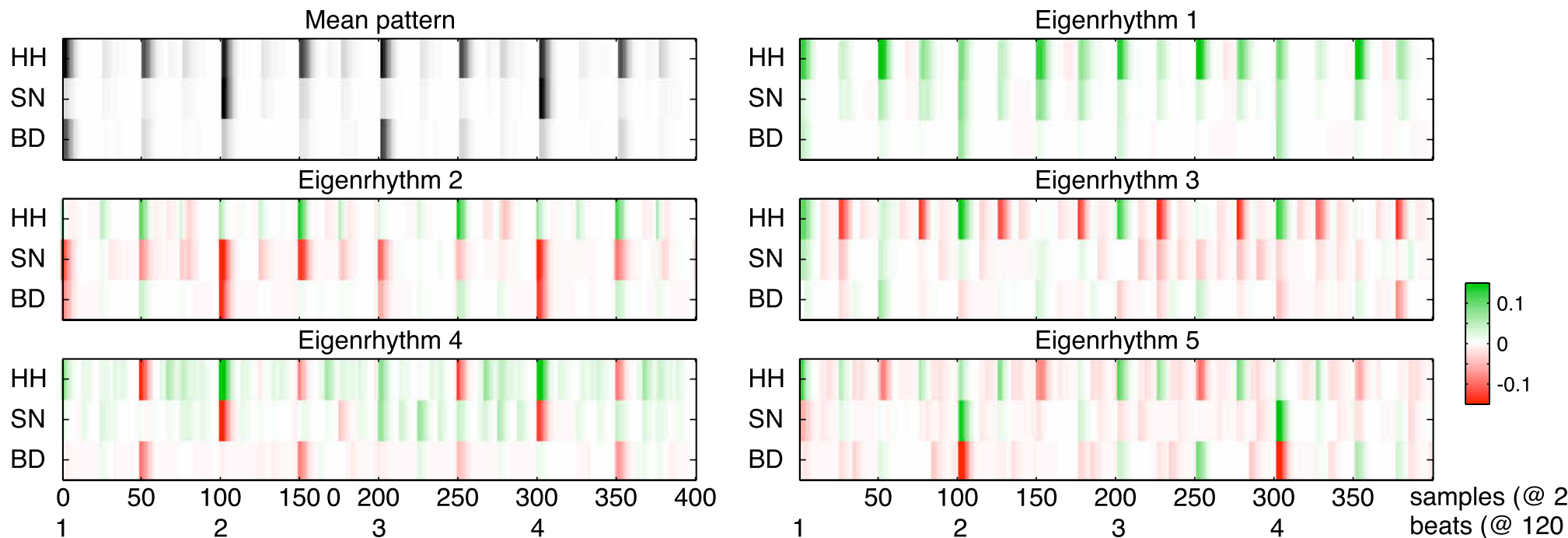
tempo (stretch):
by inferring BPM &
normalizing

downbeat (shift):
correlate against
'mean' template



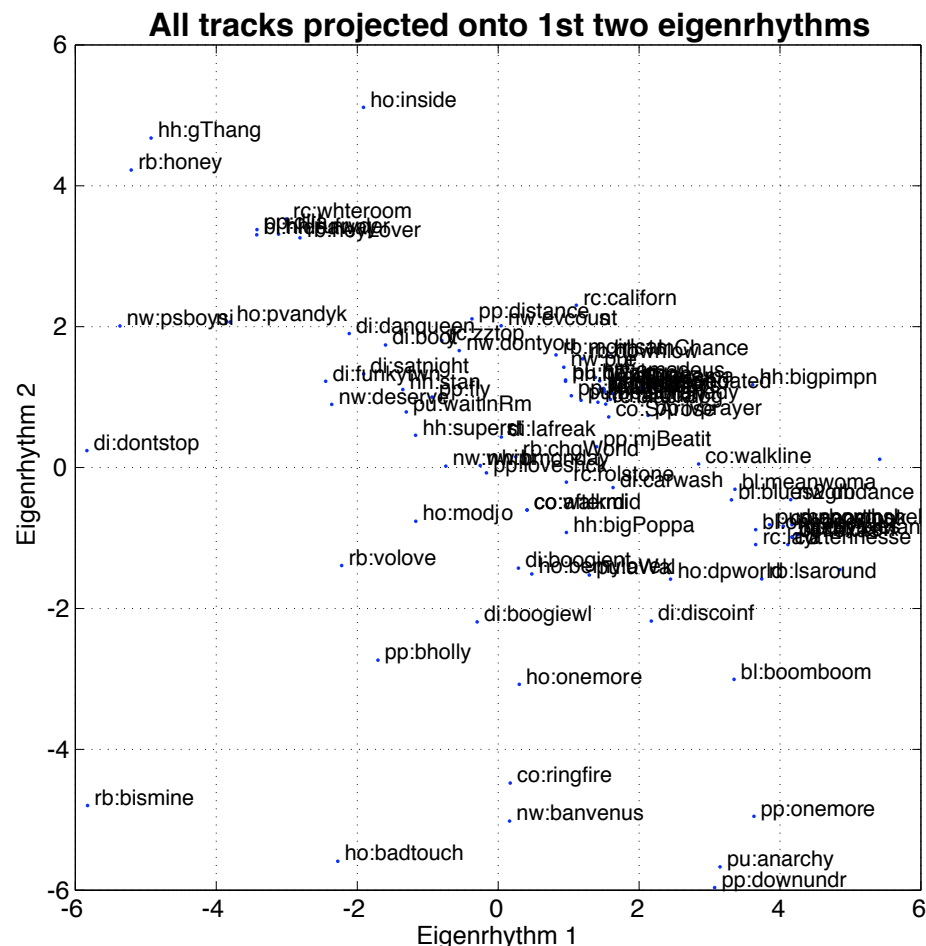
Eigenrhythms

- Need 20+ Eigenvectors for good coverage of 100 training patterns (1200 dims)
- Top patterns:



Eigenrhythms for Classification

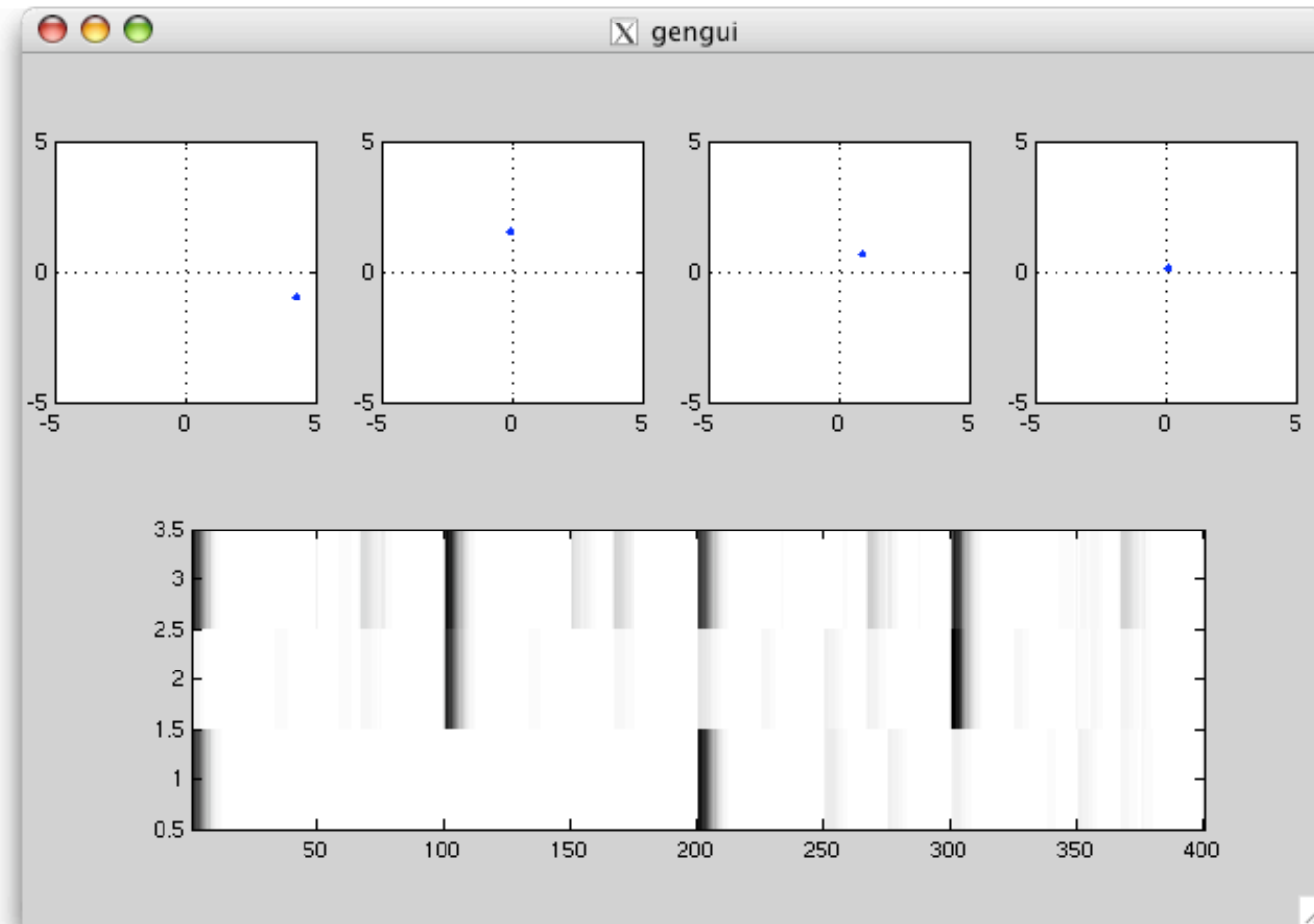
- Clusters in Eigenspace:



- Genre classification? (10 way)
 - nearest neighbor in 4D eigenspace: 21% correct

Eigenrhythm BeatBox

- Resynthesize rhythms from eigen space



2. Frequency-Domain Lin. Pred.

with Marios Athineos

- (Time-domain) Linear Prediction
 - the well-known spectral estimator

$$\rightarrow \boxed{\begin{array}{c} \text{TDLP} \\ y[n] = \sum_{i=1..p} a_i y[n-i] + e[n] \end{array}} \rightarrow$$

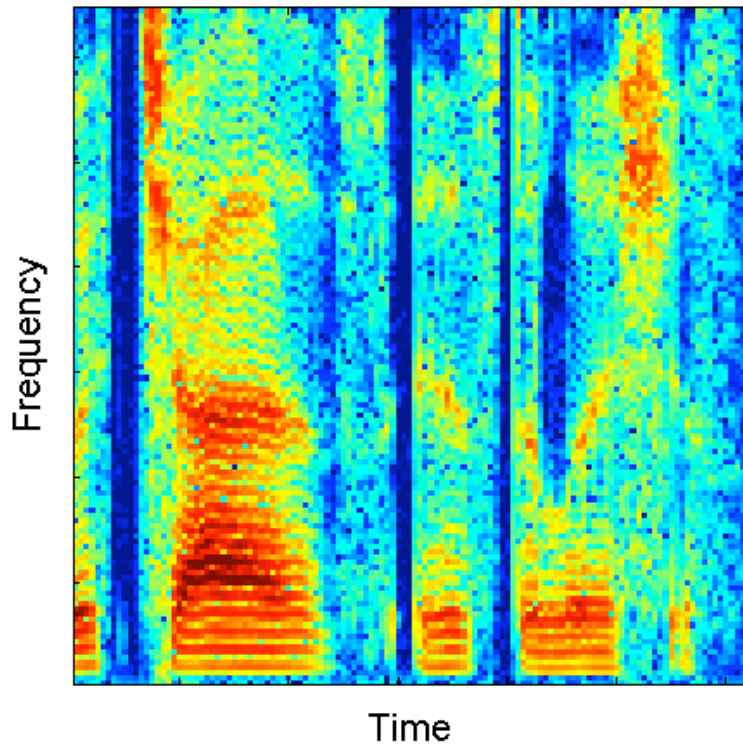
- Apply to a 'frequency domain' signal
 - dual: estimates temporal envelope

$$\rightarrow \boxed{\text{DCT}} \rightarrow \boxed{\begin{array}{c} \text{FDLP} \\ Y[k] = \sum_{i=1..p} b_i Y[k-i] + E[k] \end{array}} \rightarrow$$

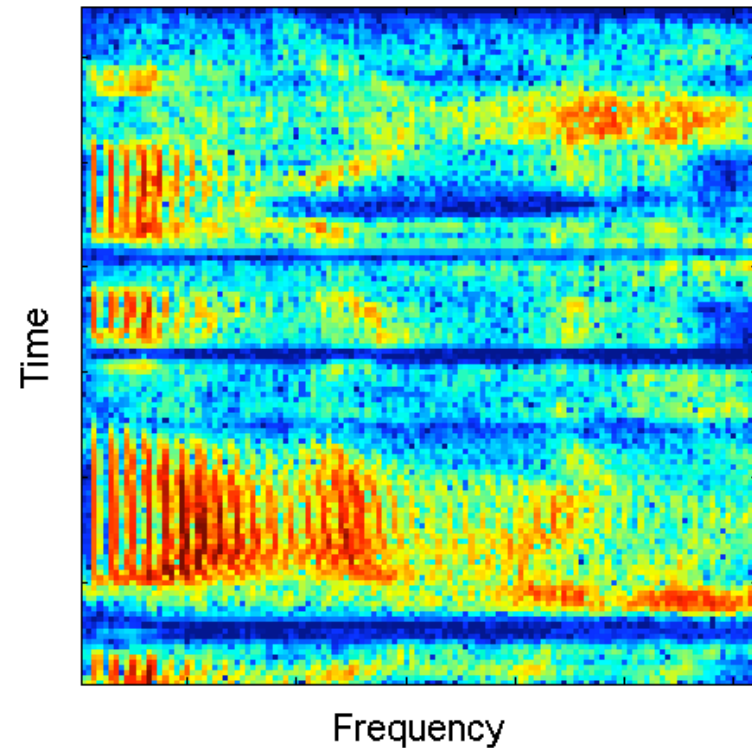
Aside: Spectrogram of the DCT

- DCT gives a pure-real signal:
Can we treat it like a waveform?

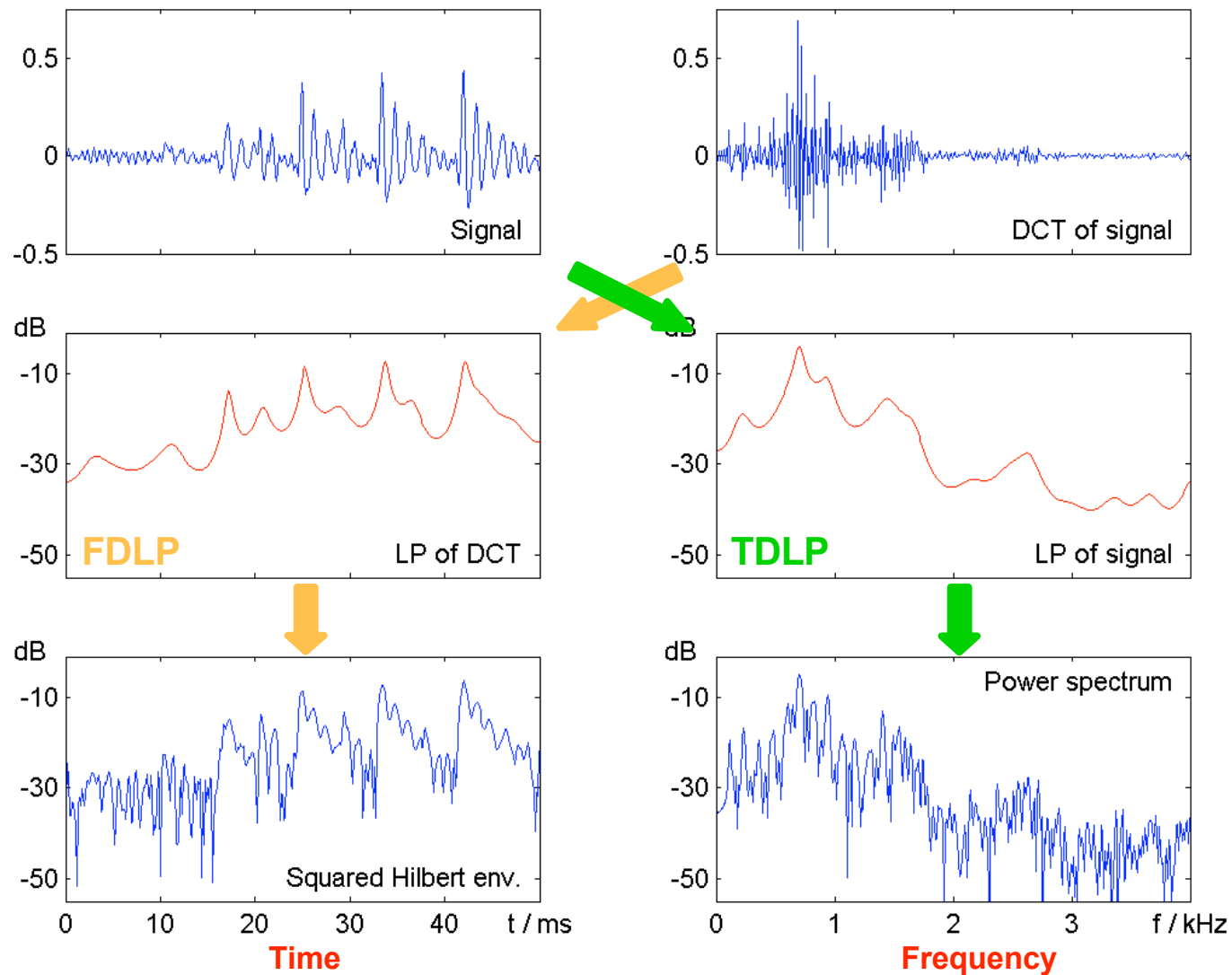
Spectrogram



DCT Spectrogram



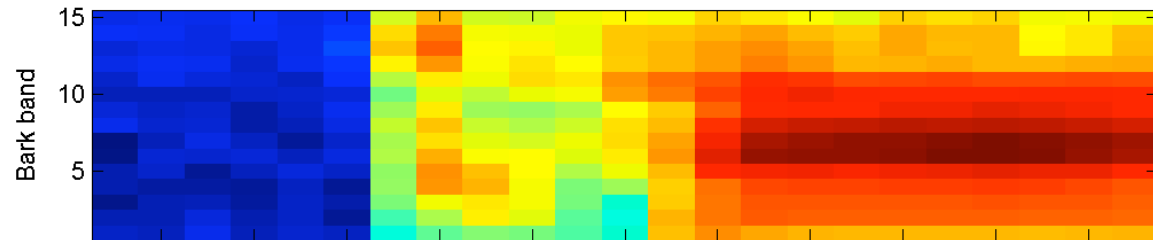
FDLP and TDLP Duality



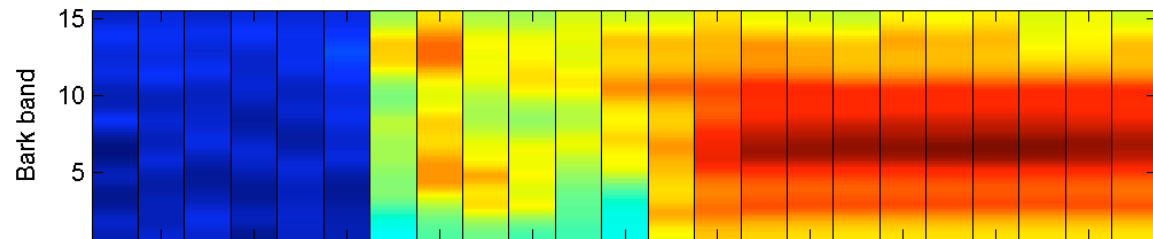
Subband FDLP

- Temporal envelopes without 25 ms windows

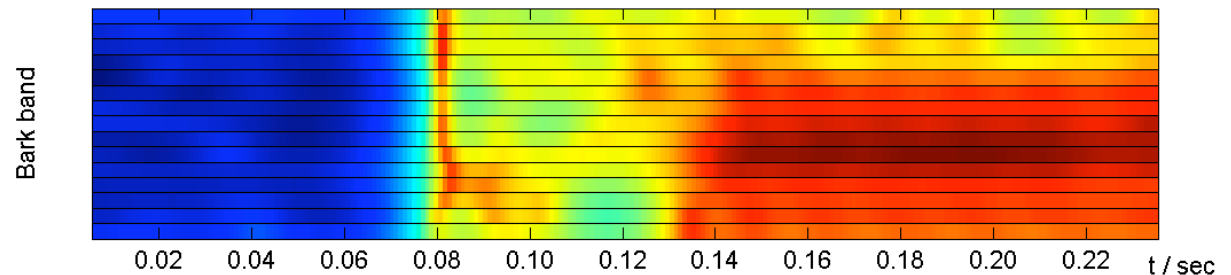
Auditory STFT
(10-25ms + Bark bin)



TDLP
(per time frame)

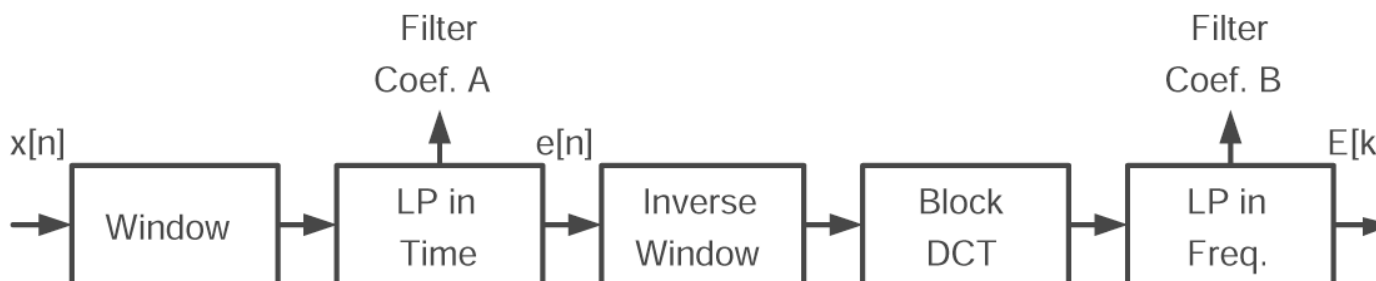


Subband FDLP
(per frequency subband)

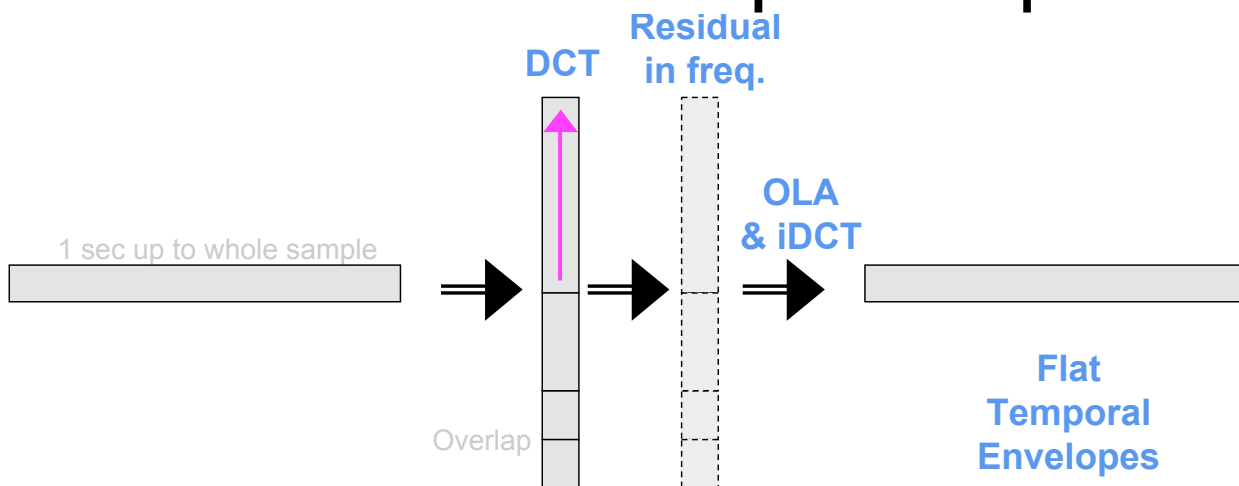


FDLP Applications

- Time-scale modification



- Modulation-domain “temporal equalization”

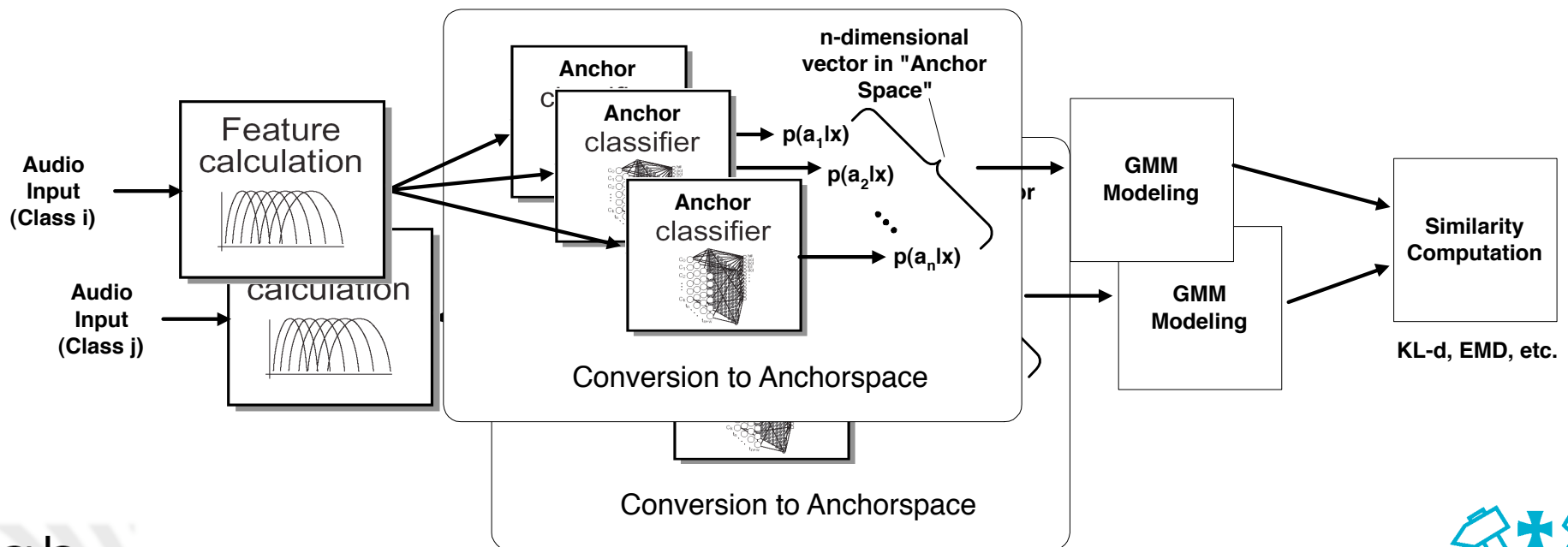


- Perceptual audio features... “PLP-squared”

3. Music Similarity Browsing

with Adam Berenzweig

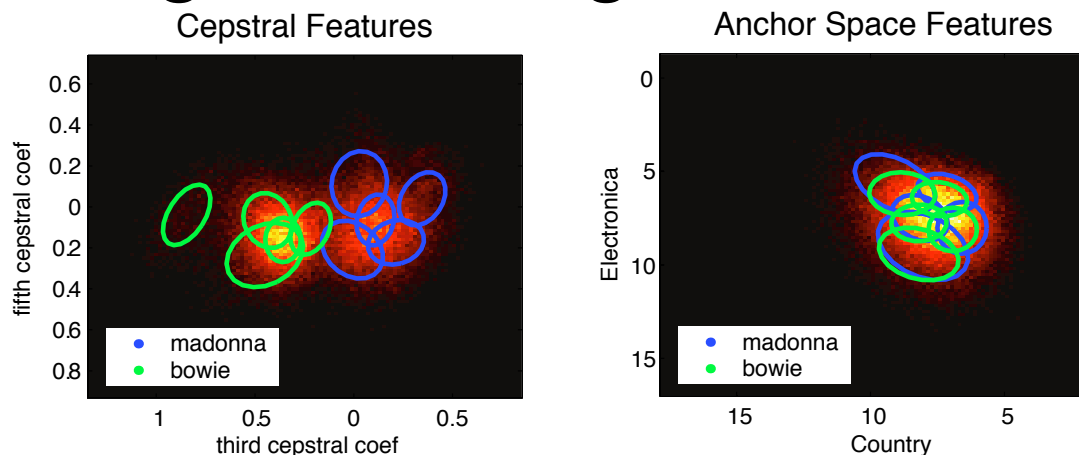
- **Musical information overload**
 - record companies filter/categorize music
 - an automatic system would be less odious
- **Connecting audio and preference**
 - map to a 'semantic space'?



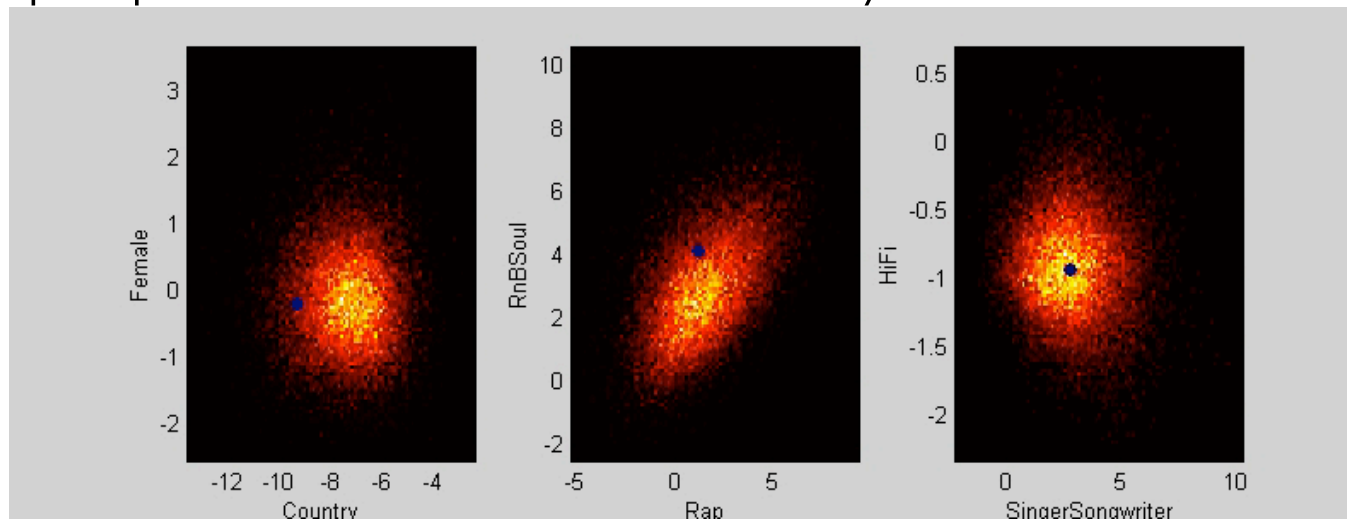
Anchor Space

- Frame-by-frame high-level categorizations

- compare to raw features?



- properties in distributions? dynamics?



'Playola' Similarity Browser

Playola Search: Artist [About] [Help] [Turn Samples Off] [Turn Debug On] [Turn Popups Off] [Logout dpwe]

Get Playola Selections: 20 songs you recently heard Go! 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	The Woodbury Muffin Outbreak	4:00	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
<input type="checkbox"/>	Monkey Dreams	The Woodbury Muffin Outbreak	2:57	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
<input type="checkbox"/>	A Cold Dark Night (Live)	The Woodbury Muffin Outbreak	3:13	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
<input type="checkbox"/>	Leo, The Ballad of	The Woodbury Muffin Outbreak	1:48	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
<input type="checkbox"/>	Baby I Forgot To Tell You	The Woodbury Muffin Outbreak	4:04	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>

Music-Space Browser [What's This?]

Feature	Less	More
AltNGrunge	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
CollegeRock	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
Country	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
DanceRock	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
Electronica	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
MetalNPunk	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
NewWave	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
Rap	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
RnBSoul	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
SingerSongwriter	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
SoftRock	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
TradRock	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
Female	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>
HiFi	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>	<div style="width: 100%; height: 10px; background-color: #ccc;"></div>

Similar Songs: [Play this list] [What's This?]

	Song Title	Artist	Distance	Good Match?
<input type="checkbox"/>	Baby I Forgot To Tell You	The Woodbury Muffin Outbreak	0.00	
<input type="checkbox"/>	Number five	Bizi Chyld	0.07	
<input type="checkbox"/>	Waiting for Your Love	Toto	0.08	
<input type="checkbox"/>	Excerpt from 'CD'	Weirdomusic	0.08	



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](#)

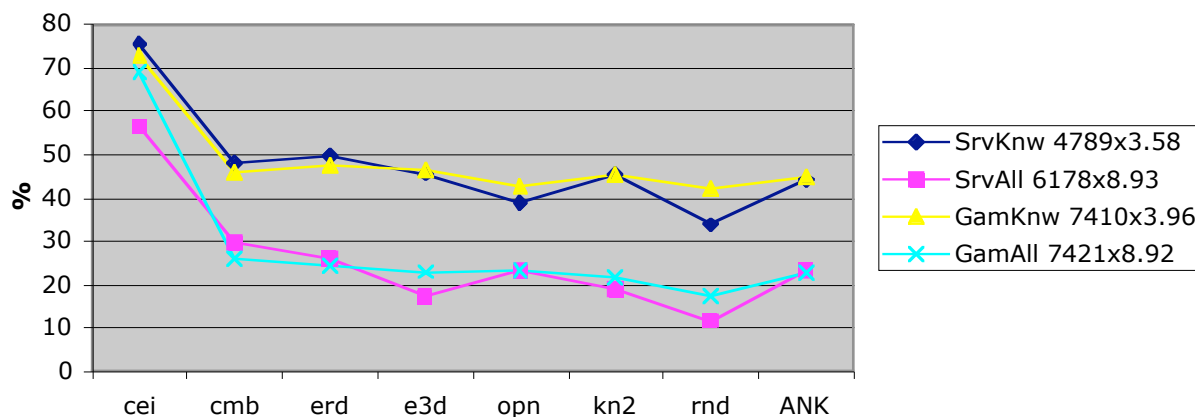
Evaluation

- Compare Anchor Space measures against Musicseer subjective results

- “triplet” agreement percentage
- Top-N ranking agreement score:

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

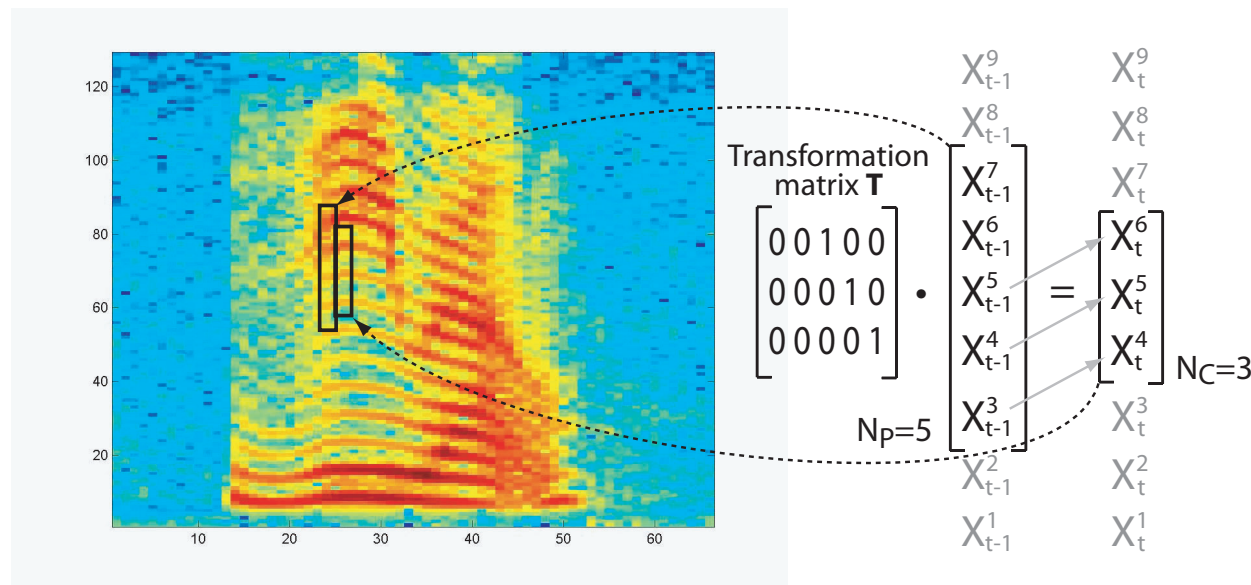
- First-place agreement percentage
 - simple significance test



4. Transformation-based models

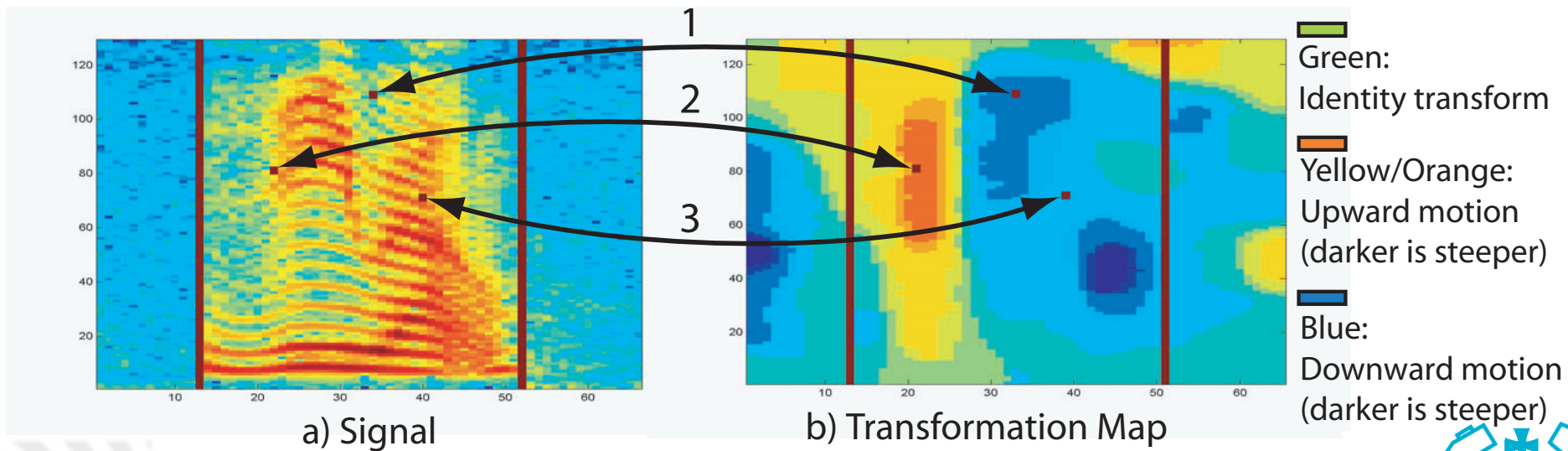
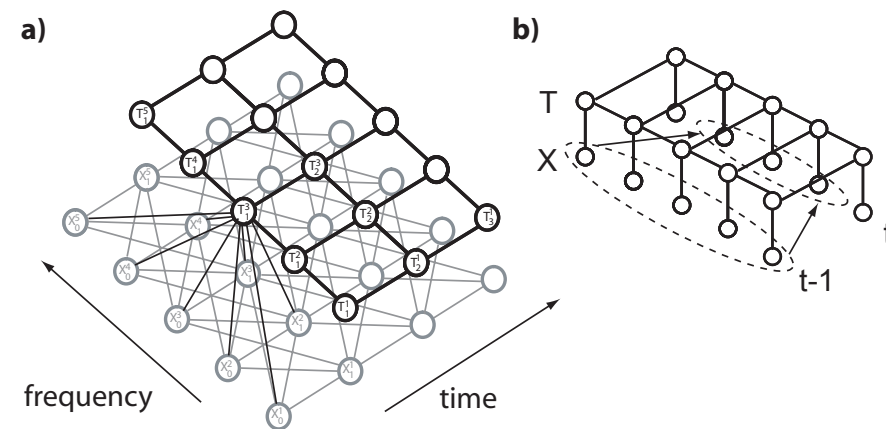
with Manuel Reyes and Nebojsa Jojic

- HMMs are poor generative models
 - accurate modeling requires 1000s of states
- Observation:
 - Speech spectra undergo minor **deformations**
 - suggests a different generative model?



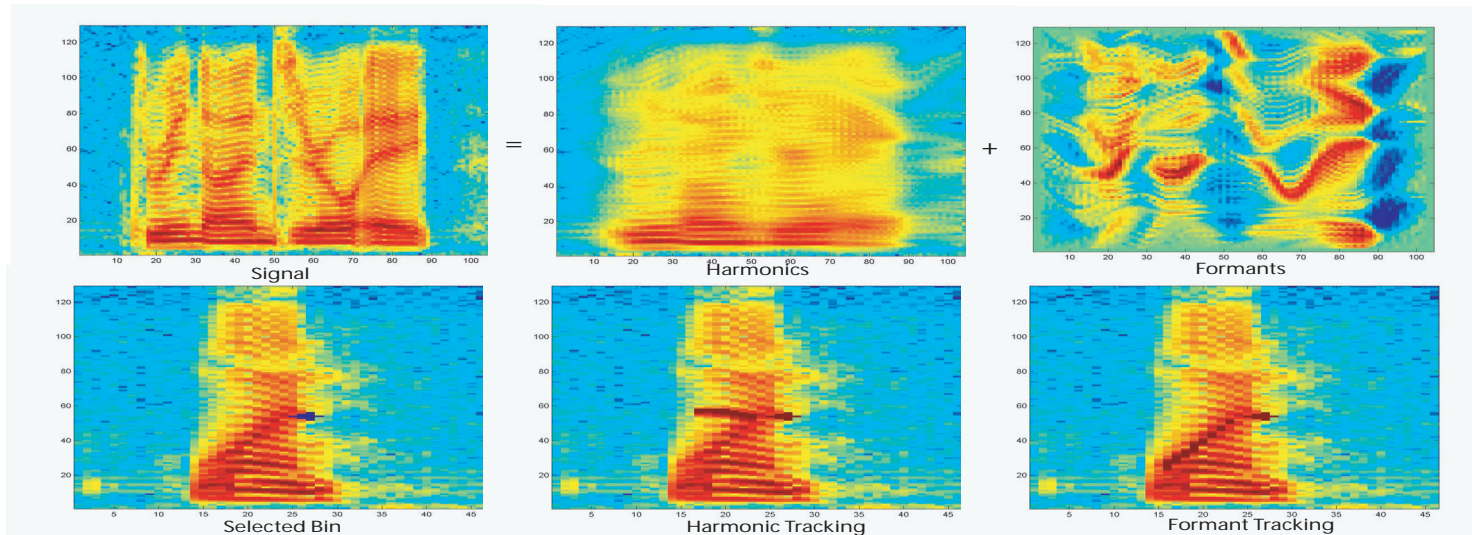
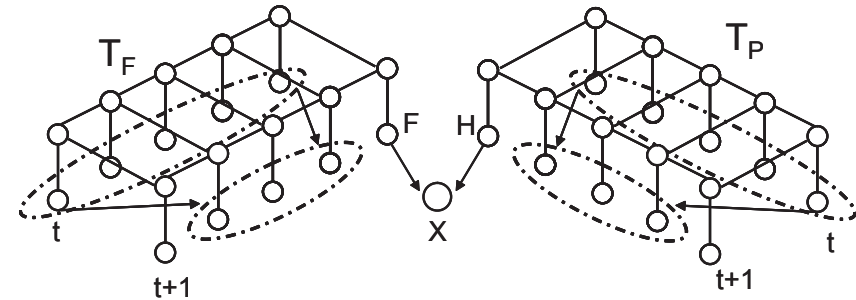
States+Transformation Model

- Time-frequency state grid
- State \rightarrow explicit prototype or a transformation on prior frame
- Infer underlying states



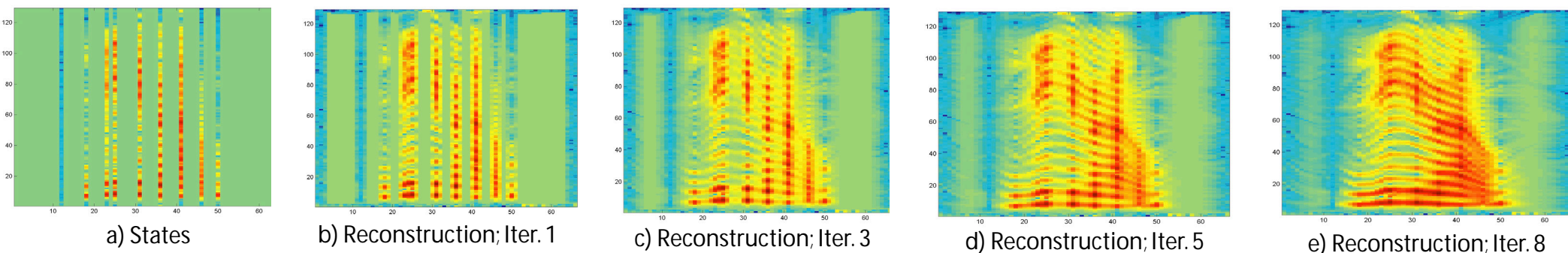
Two-layer model

- Source-filter decomposition
 - pitch and formants have different dynamics
- Apply transformation models for both
 - log-spectra:
sum of excitation & filter
 - inference does separation

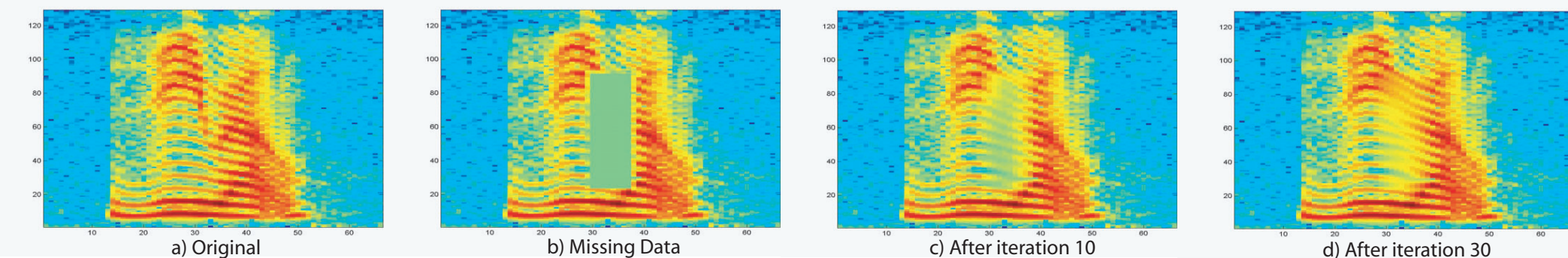


Transformation model applications

- Compact, accurate source descriptions
 - only a few explicit states needed



- Belief propagation can infer missing values
 - .. of state grid, hence magnitude spectrum



5. Segmenting Personal Audio

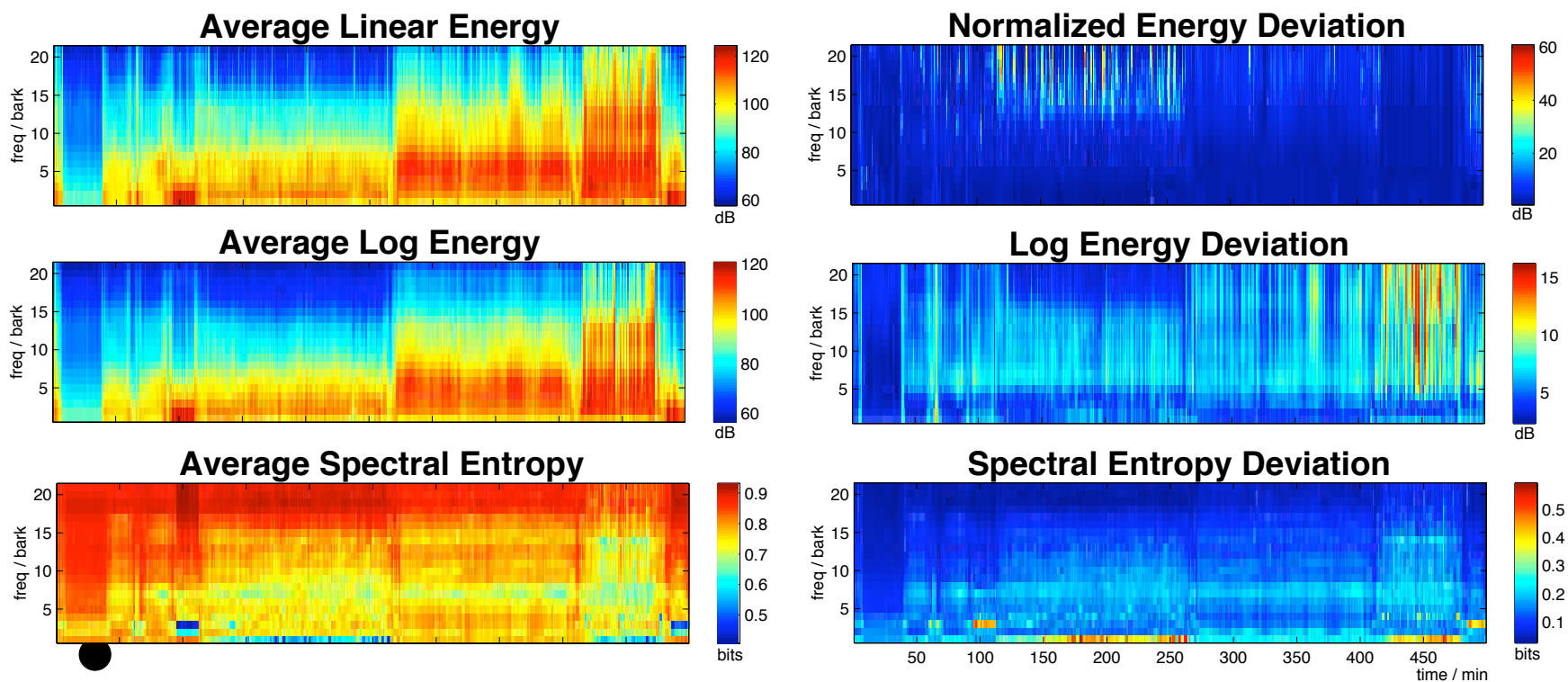
with Kean sub Lee

- Easy to record **everything** you hear
 - ~100GB / year @ 64 kbps
- Very hard to **find anything**
 - how to scan?
 - how to visualize?
 - how to index?
- Starting point: Collect **data**
 - ~ 60 hours (8 days, ~7.5 hr/day)
 - hand-mark 139 segments (26 min/seg avg.)
 - assign to 41 classes (8 have multiple instances)



Features for Long Recordings

- Feature frames = 1 min (not 25 ms!)
- Characterize variation within each frame...



○ and structure within coarse auditory bands

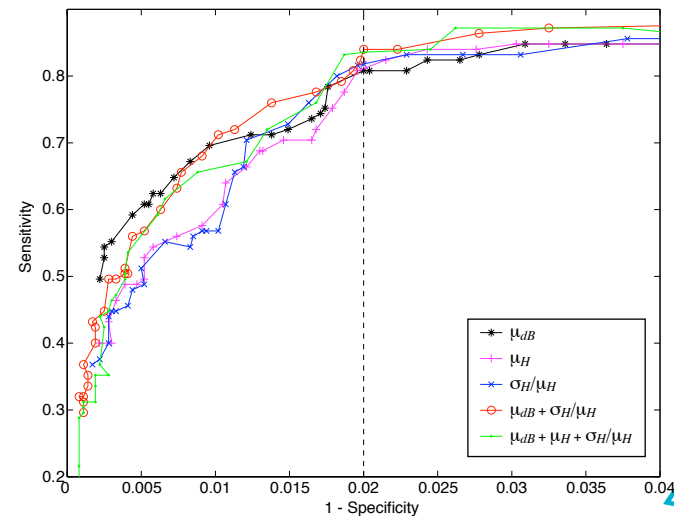
BIC Segmentation

- **Untrained segmentation technique**
 - statistical test indicates good change points:

$$\log \frac{L(X_1; M_1)L(X_2; M_2)}{L(X; M_0)} \geq \frac{\lambda}{2} \log(N) \Delta \#(M)$$

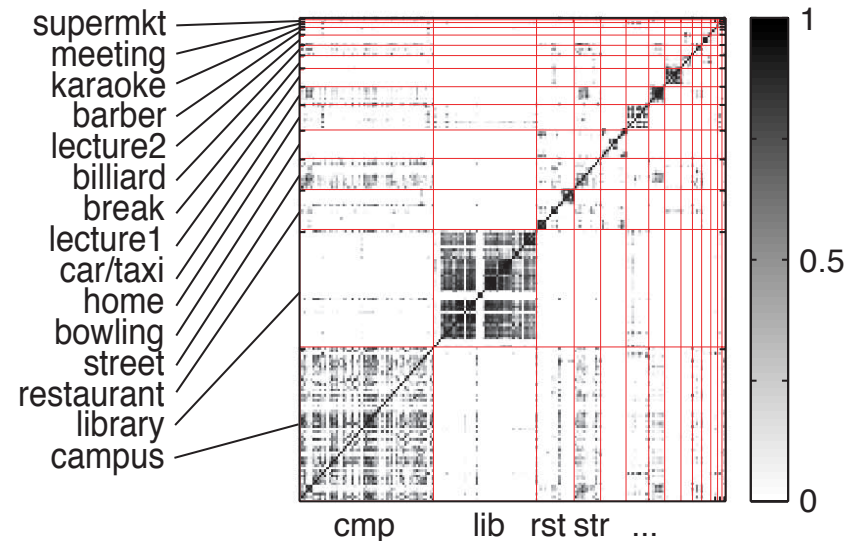
- **Evaluate: 60hr hand-marked boundaries**
 - different features & combinations
 - Correct Accept % @ False Accept = 2%:

μ_{dB}	80.8%
μ_H	81.1%
σ_H/μ_H	81.6%
$\mu_{dB} + \sigma_H/\mu_H$	84.0%
$\mu_{dB} + \sigma_H/\mu_H + \mu_H$	83.6%



Segment clustering

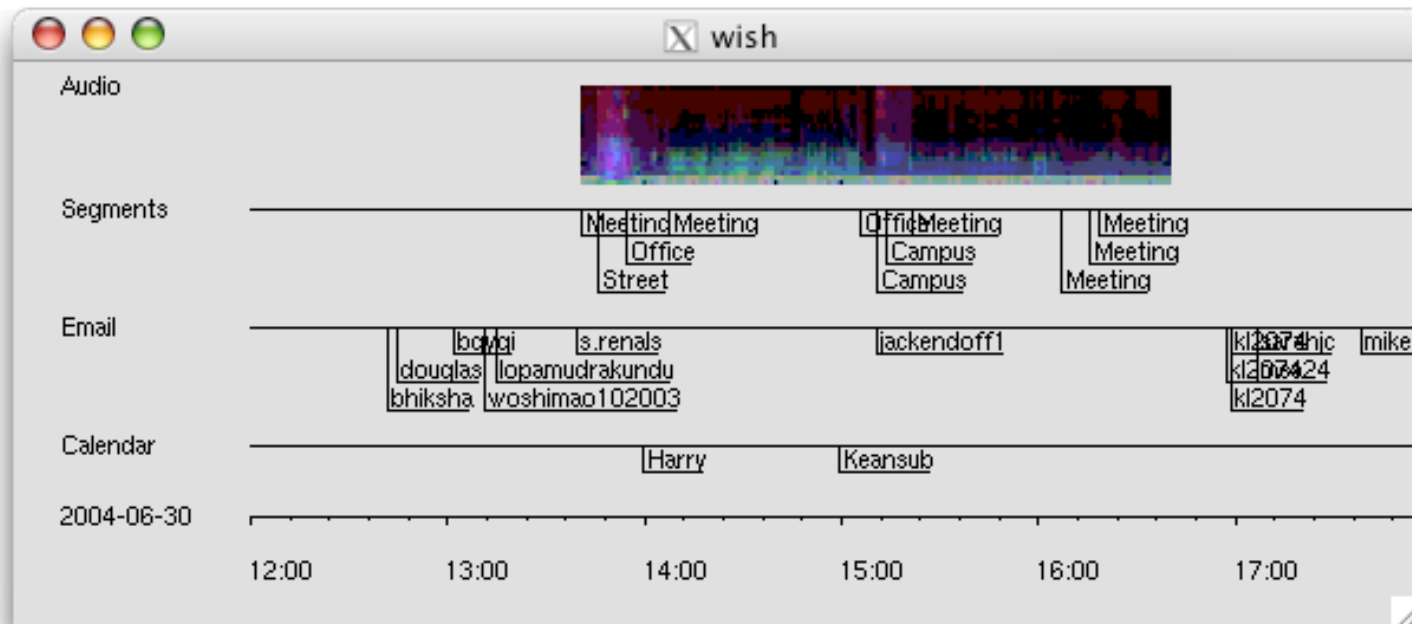
- Daily activity has lots of repetition:
Automatically cluster similar segments



- Spectral clustering achieves ~60% correct
 - 16-way ground truth labels

Future Work

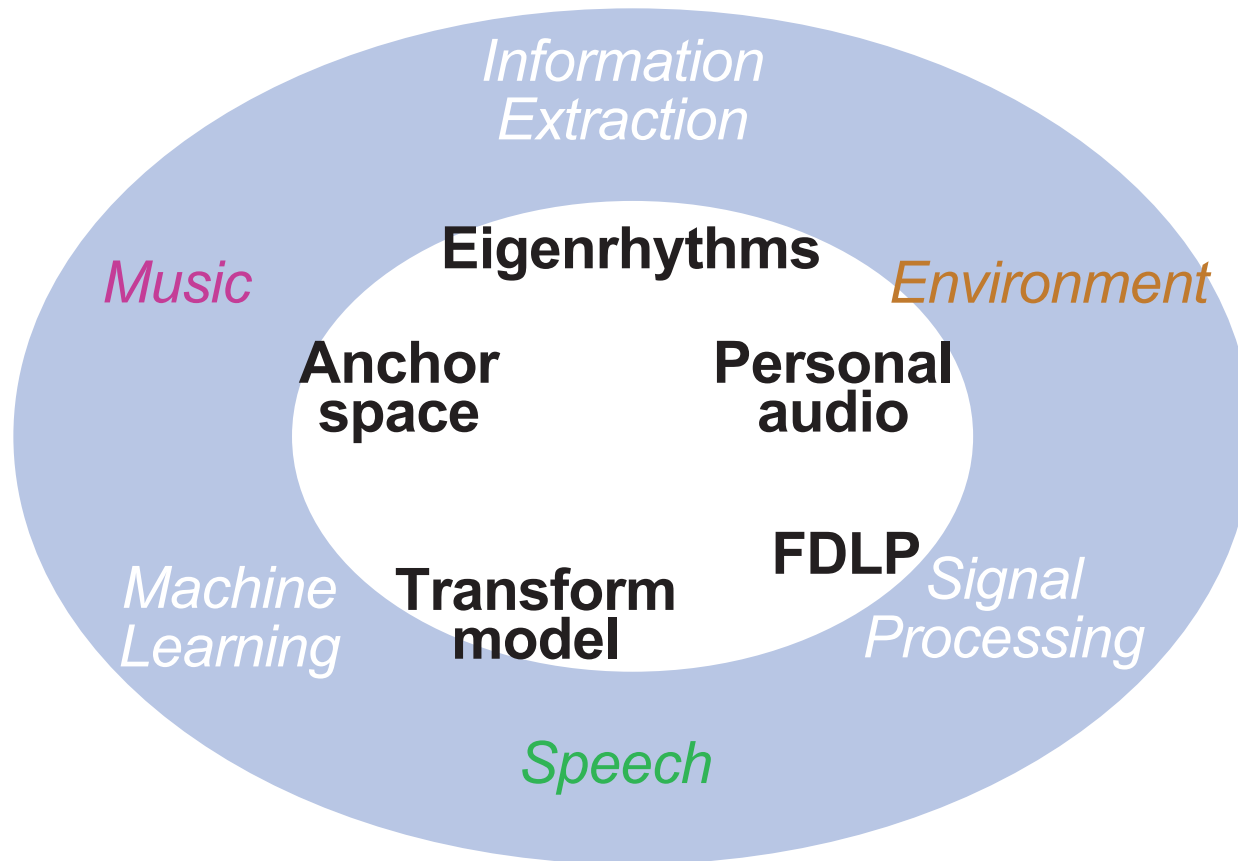
- **Visualization** / browsing / diary inference
 - link to other information sources



- **Privacy protection**
 - speaker/speech “**search and destroy**”

Summary

- Today's topics:



- + Speech recognition, Meeting recordings

Audio/Music @ LabROSA - Dan Ellis

2004-07-29

LabROSA Summary

- **LabROSA**
 - signal processing
 - + machine learning
 - + info extraction
- **Applications**
 - Eigenrhythms: drum pattern models
 - FDLP temporal envelope models
 - Music Similarity Browsing
 - Transformation-based generative models
 - Personal audio analysis
- **Also...**
 - speech recognition, meeting recordings, ...

