

DAY 3

Intelligent Audio Systems: A review of the foundations and applications of semantic audio analysis and music information retrieval



Jay LeBoeuf
Imagine Research
jay@imagine-research.com

Kyogu Lee
Gracenote
Klee@ccrma.stanford.edu

June 2009

These lecture notes contain hyperlinks to the CCRMA Wiki.

On these pages, you can find supplemental material for lectures - providing extra tutorials, support, references for further reading, or demonstration code snippets for those interested in a given topic .

Click on the  symbol on the lower-left corner of a slide to access additional resources.

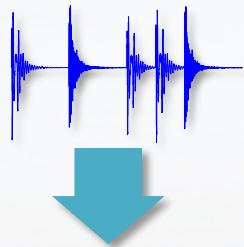
WIKI REFERENCES...



Review from Day 2

- BBQ Today
 - Correction on knn formatting
 - Name some spectral features
 - What are the 3 major components of a MIR system?
 - Why do we have to scale our extracted features?
 - Which of these did we really not do at all in Lab 2? And, do you think this was a problem?
-
- How did the lab go?
 - Let's dig into some interesting observations from the lab
 - Did you try other audio files – other instrument recognizers?

Basic system overview



Segmentation

(Frames, Onsets,
Beats, Bars, Chord
Changes, etc)



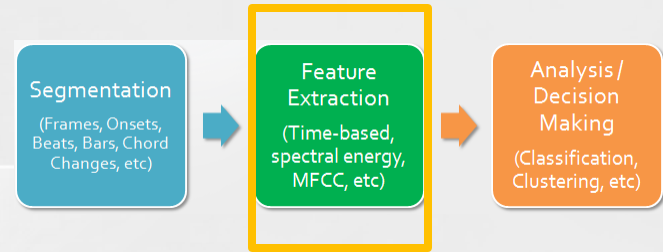
Feature Extraction

(Time-based,
spectral energy,
MFCC, etc)



Analysis / Decision Making

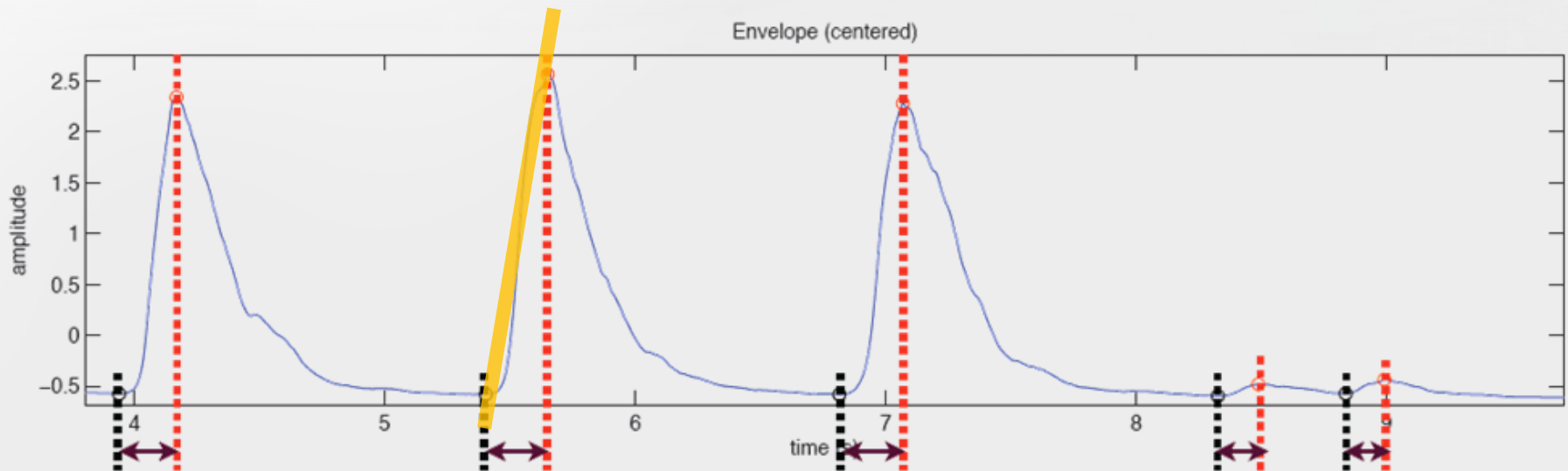
(Classification,
Clustering, etc)



FEATURE EXTRACTION

Temporal Information

- Rise time or Attack time- time interval between the onset and instant of maximal amplitude
- Attack slope

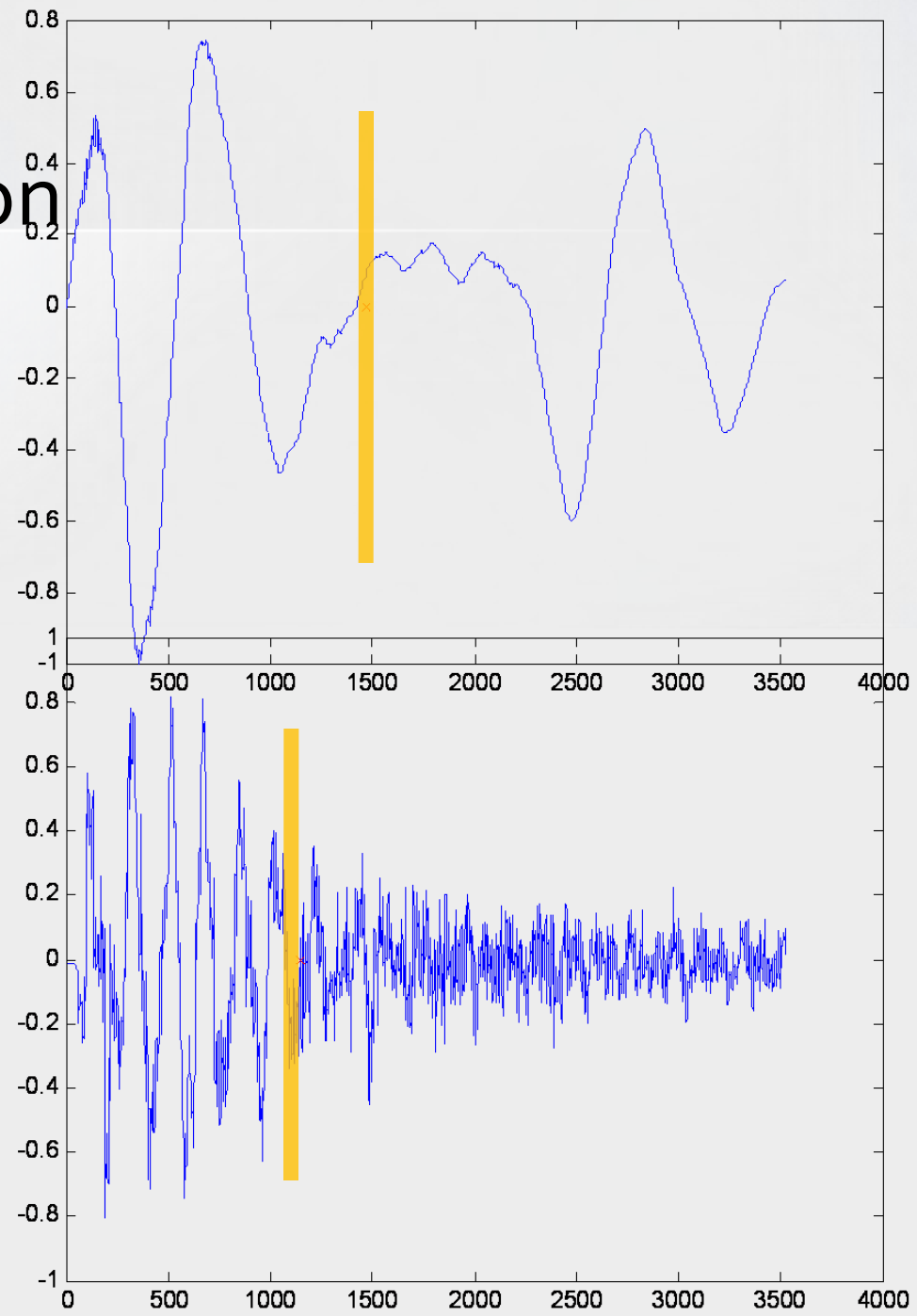


Picture courtesy: Olivier Lartillot



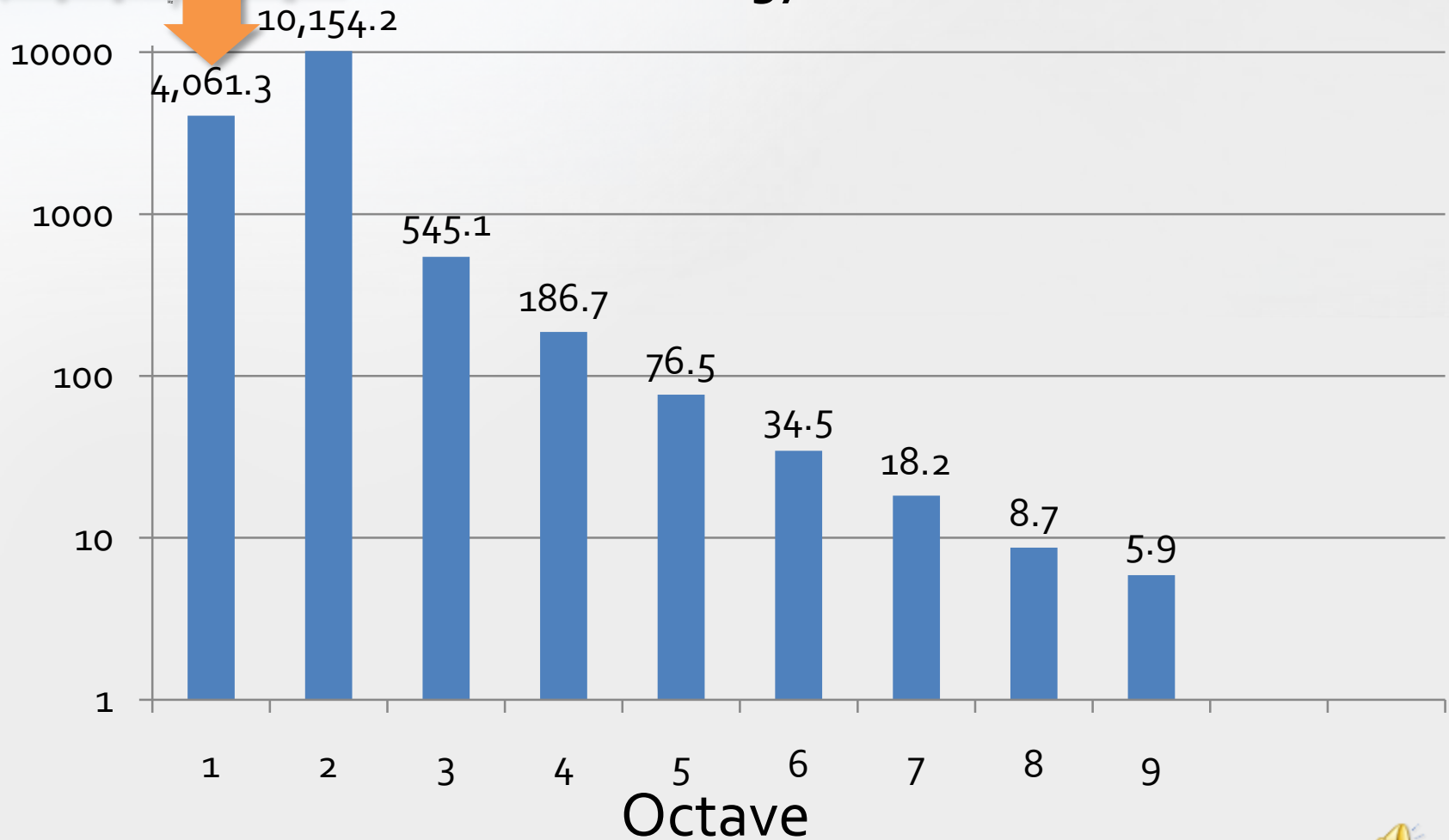
Temporal Information

- Temporal Centroid



Frame 1

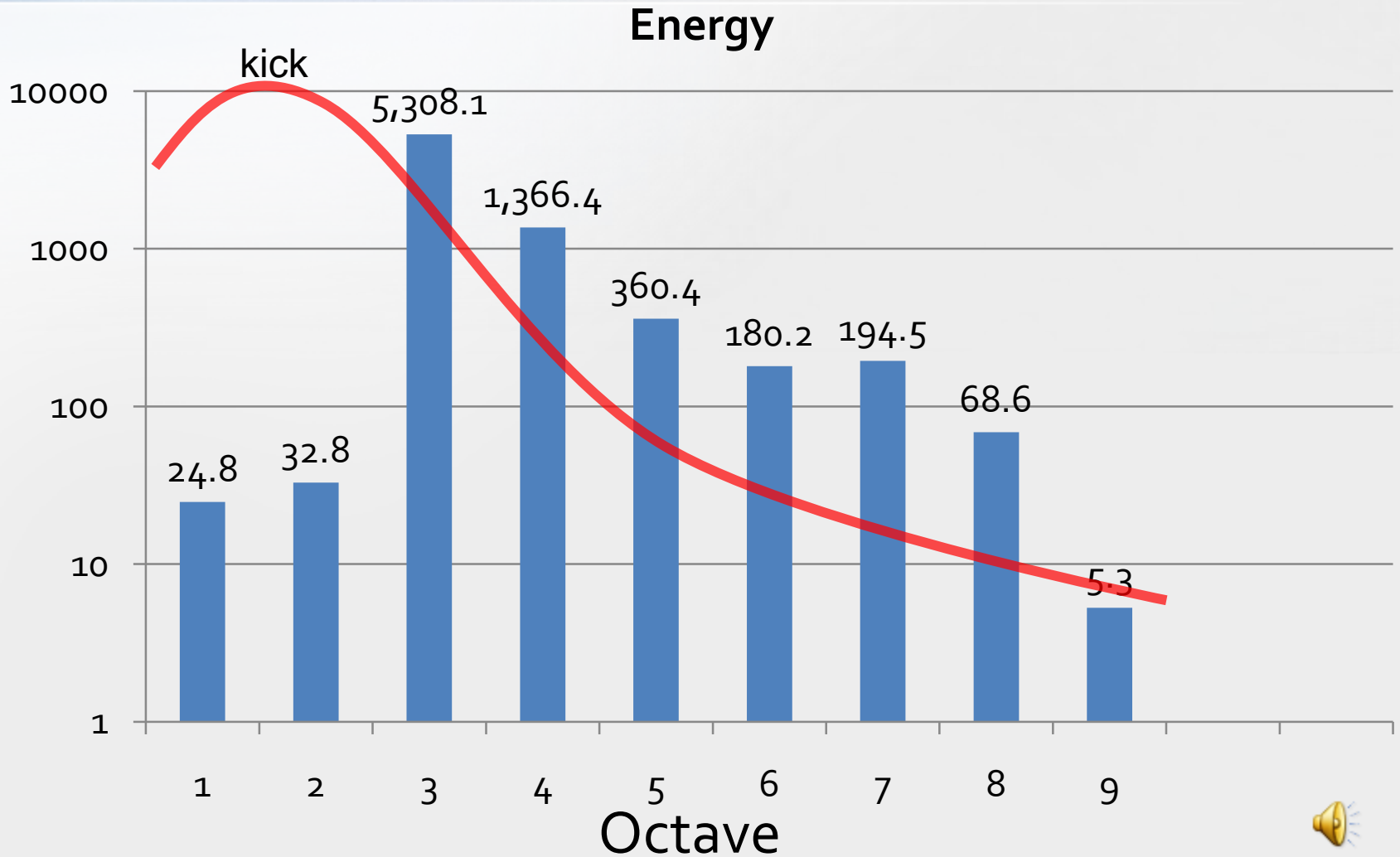
Energy



Features – Frame 1

Frame	ZC R	Centroid	BW	Skew	Kurtosis	E1	E2	E3	E4	E5	E6	E7	E8	E9
1	9	2.8kHz	5kHz	2.2	6.7	4000	10100	545	187	77	35	18	9	6

Frame 2



Features : SimpleLoop.wav

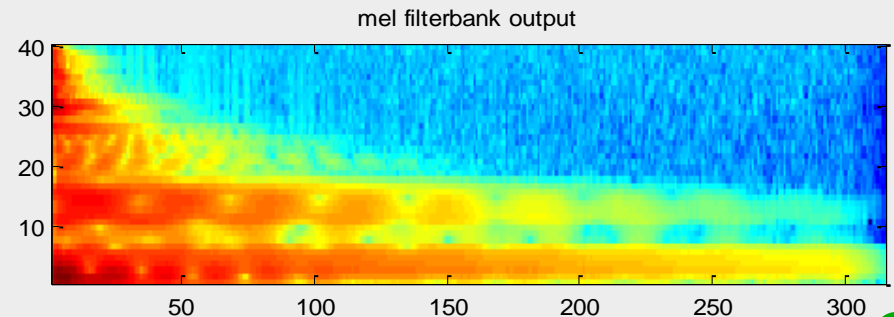
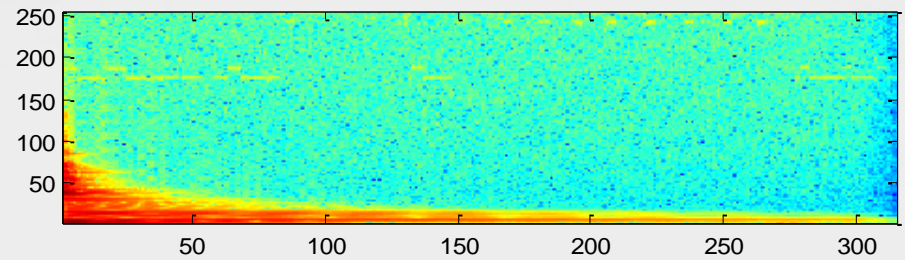
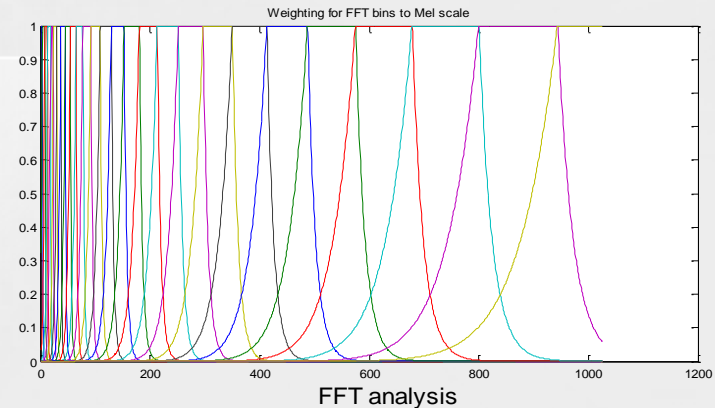
Frame	ZC R	Centroid	BW	Skew	Kurtosis	E1	E2	E3	E4	E5	E6	E7	E8
1	9	2.8kHz	5kHz	2.2	6.7	4000	10100	545	187	77	35	18	9
2	423	3.1kHz	4kHz	2	7.2	24	33	5300	1366	360	180	194	68

MFCCs

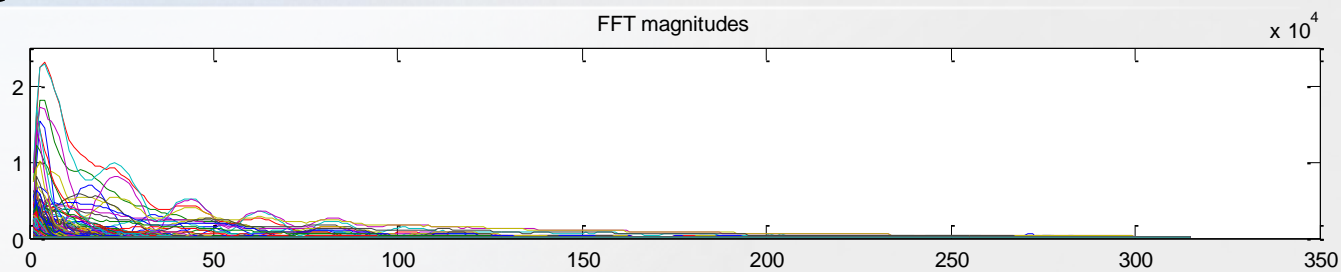
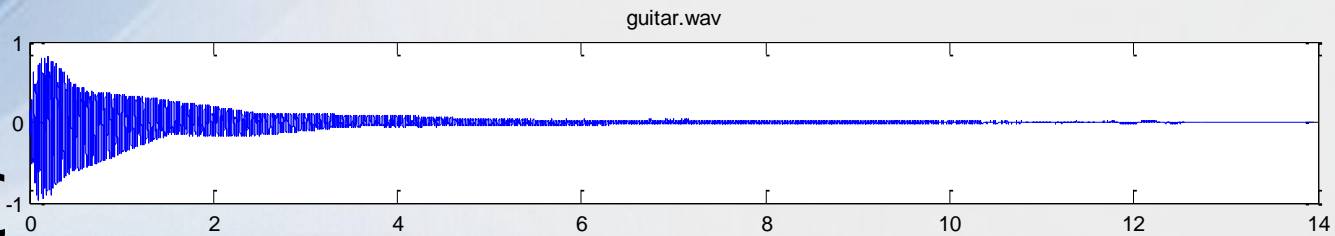
The idea of MFCCs is to capture spectrum in accordance with human perception.

1. STFT
2. $\log(\text{STFT})$
3. Perform mel-scaling to group and smooth coefficients. (perceptual weighting)
4. Decorrelate with DCT

[...continued...]



MFCC



1

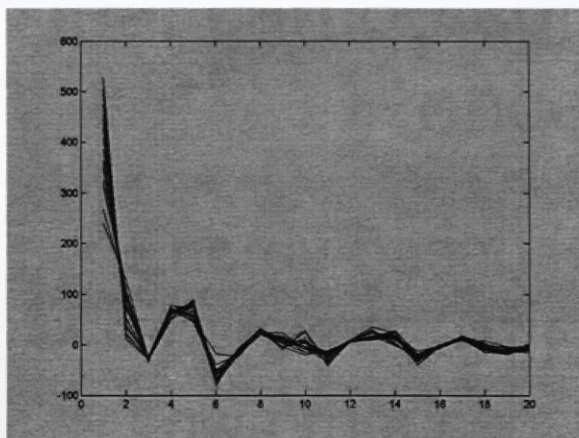
2

3

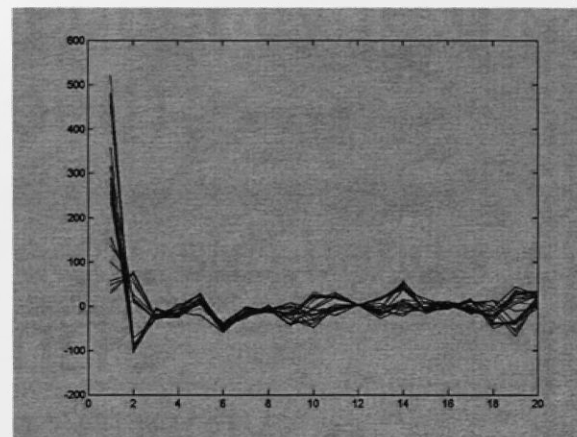
4

MFCC of Music

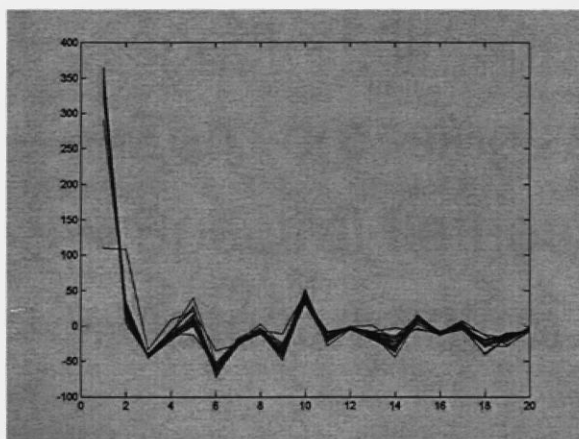
(Petruncio, 2003)



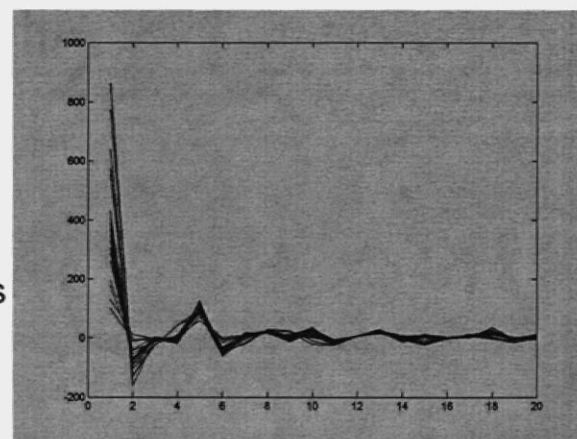
Piano



Saxophone



Tenor
Opera
Singer



Drums

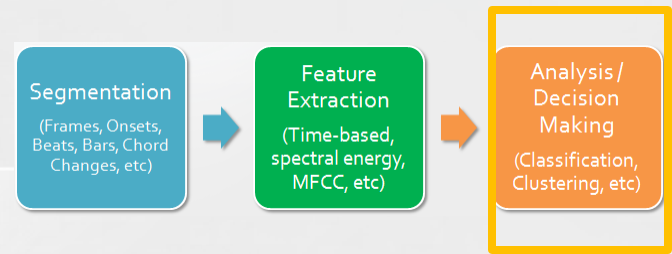
Spectral Energy vs. MFCC

Features: Measuring changes

- Δ and $\Delta \Delta$
 - Change between frames
 - How quickly the change is occurring
- Spectral flux is the distance between the spectrum of successive frames

Feature extraction

- Feature design and creation uses one's domain knowledge.
- Choosing discriminating features is critical
- Smaller feature space yields smaller, simpler models, faster training, often less training data needed



ANALYSIS AND DECISION MAKING

Supervised vs. Unsupervised

- Unsupervised - “clustering”
- Supervised – binary classifiers (2 classes)
- Multiclass is derived from binary

Clustering

- Unsupervised learning – find pockets of data to group together
- Statistical analysis techniques

Clustering

- $K = \#$ of clusters
- Choosing the number of clusters – note that choosing the “best” number of clusters according to minimizing total squared distance will always result in same $\#$ of clusters as data points.

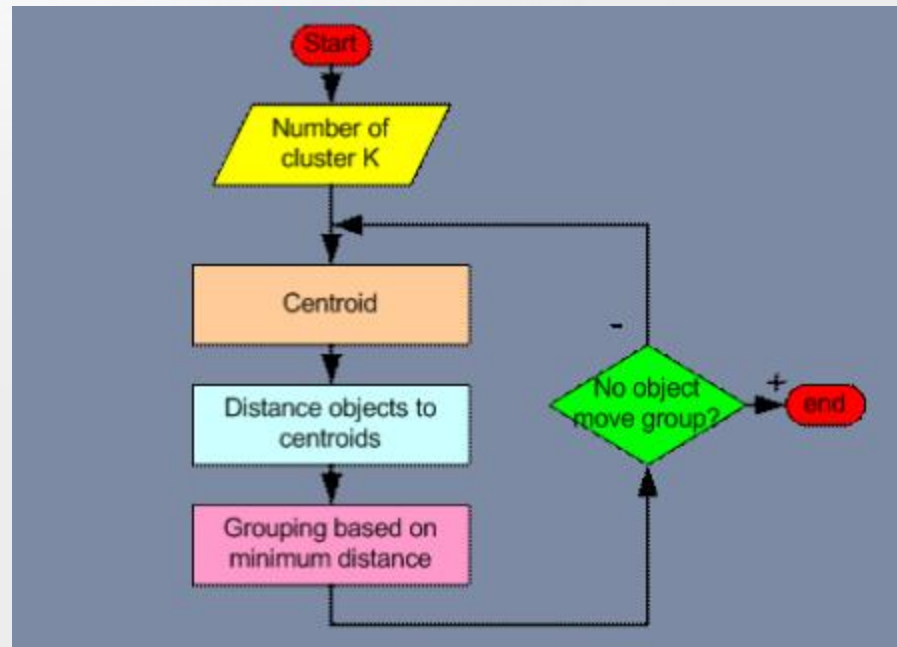
Clustering

The basic goal of clustering is to divide the data into groups such that the points within a group are close to each other, but far from items in other groups.

Hard clustering – each point is assigned to one and only one cluster.

Demo

- http://home.dei.polimi.it/matteucc/Clustering/tutorial_html/AppletKM.html



K-Means

The key points relating to *k-means clustering* are:

- k-means is an automatic procedure for clustering unlabelled data;
- it requires a pre-specified number of clusters;
- Clustering algorithm chooses a set of clusters with the minimum within-cluster variance
- Guaranteed to converge (eventually)
- Clustering solution is dependent on the initialization
(You get different results with each running)



K-Means

The initialization method needs to be further specified.
There are several possible ways to initialize the cluster centers:

- *Choose random data points as cluster centers*
- *Randomly assign data points to K clusters and compute means as initial centers*
- *Choose data points with extreme values*
- *Find the mean for the whole data set then perturb into k means*
- *Find ground-truth for data*

EVALUATION

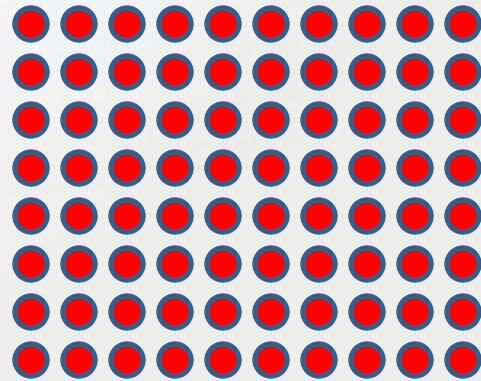
Our classifier accuracy is 83.4%

Cross-validation

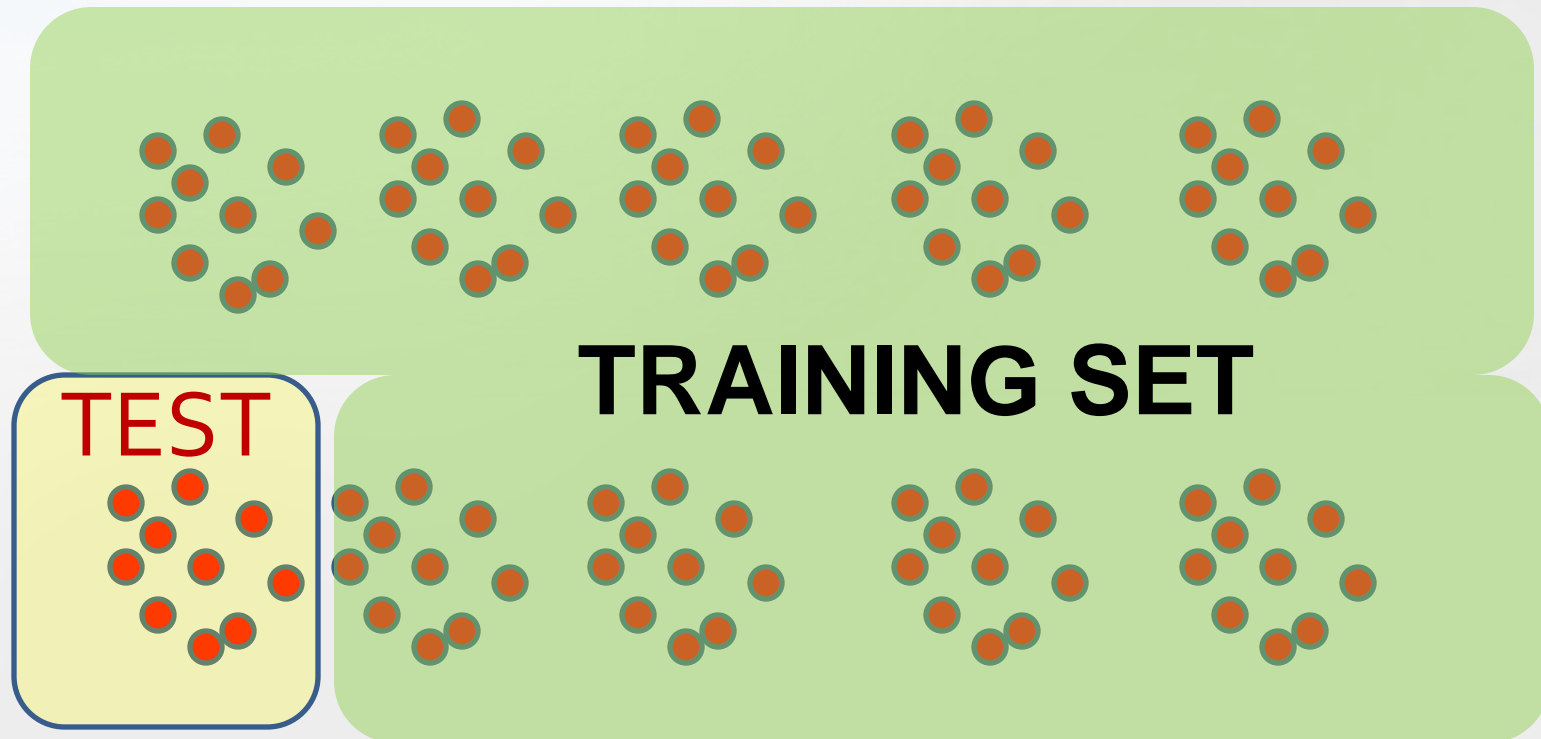
- Say, 10-fold cross validation
- Divide test set into 10 random subsets.
- 1 test set is tested using the classifier trained on the remaining 9.
- We then do test/train on all of the other sets and average the percentages. Helps prevent over fitting.
- Do not optimize too much on cross validation – you can severely overfit. Sanity check with a test set.



Cross-validation

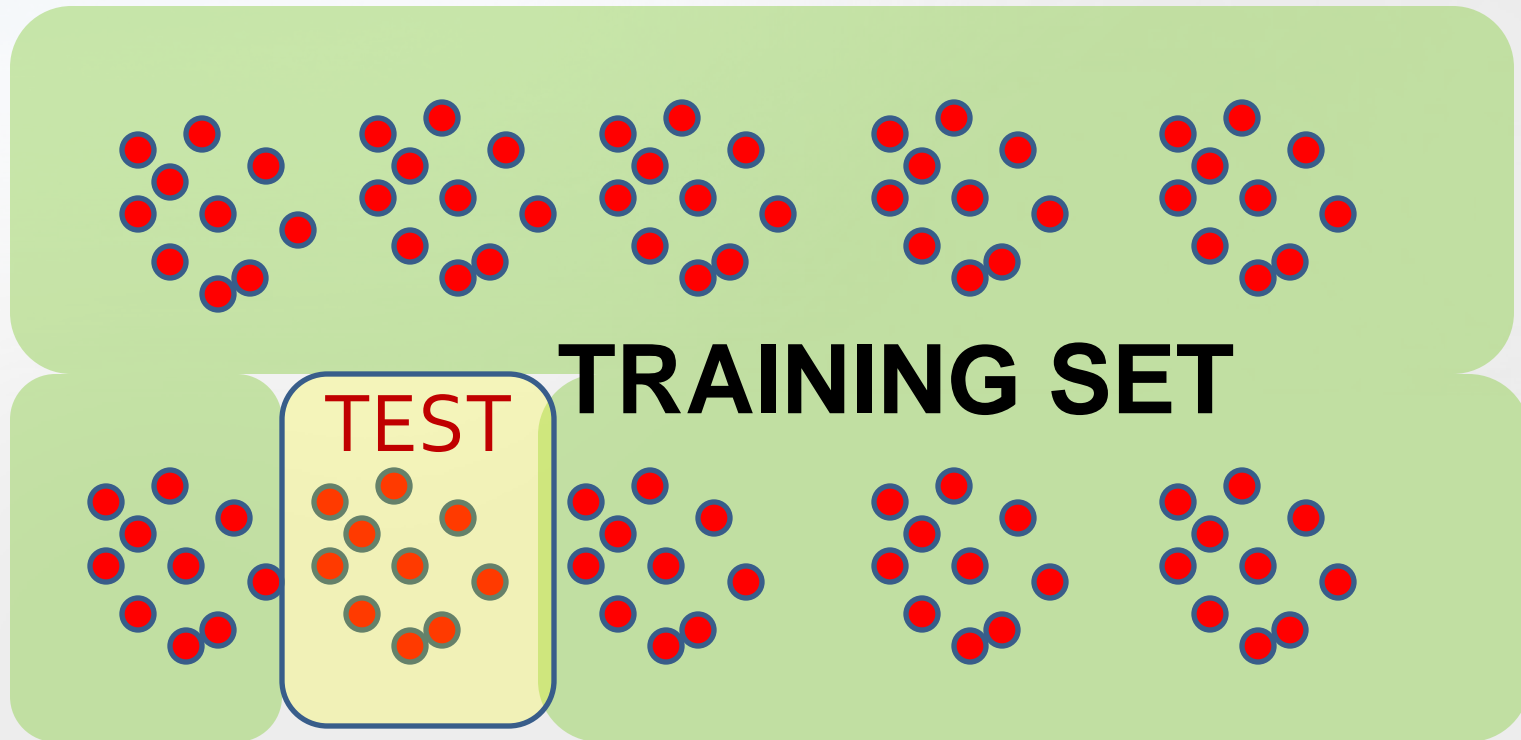


Cross-validation



Fold 1: 70%

Cross-validation



Fold 1: 70%

Fold 2: 80%

Cross-validation

Fold 1: 76%

Fold 2: 80%

Fold 3: 77%

Fold 4: 83%

Fold 5: 72%

Fold 6: 82%

Fold 7: 81%

Fold 8: 71%

Fold 9: 90%

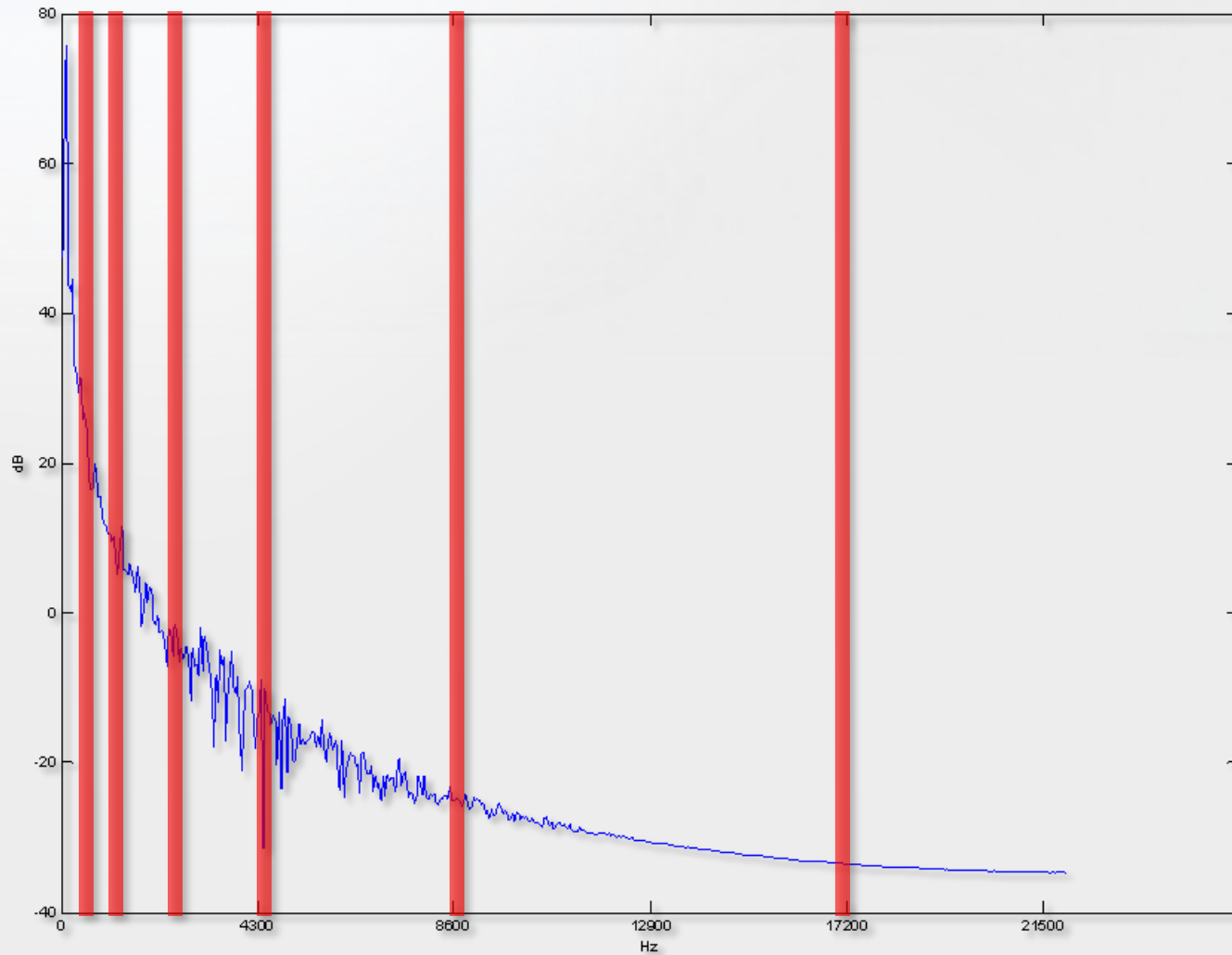
Fold 10: 82%

Mean = 79.4%

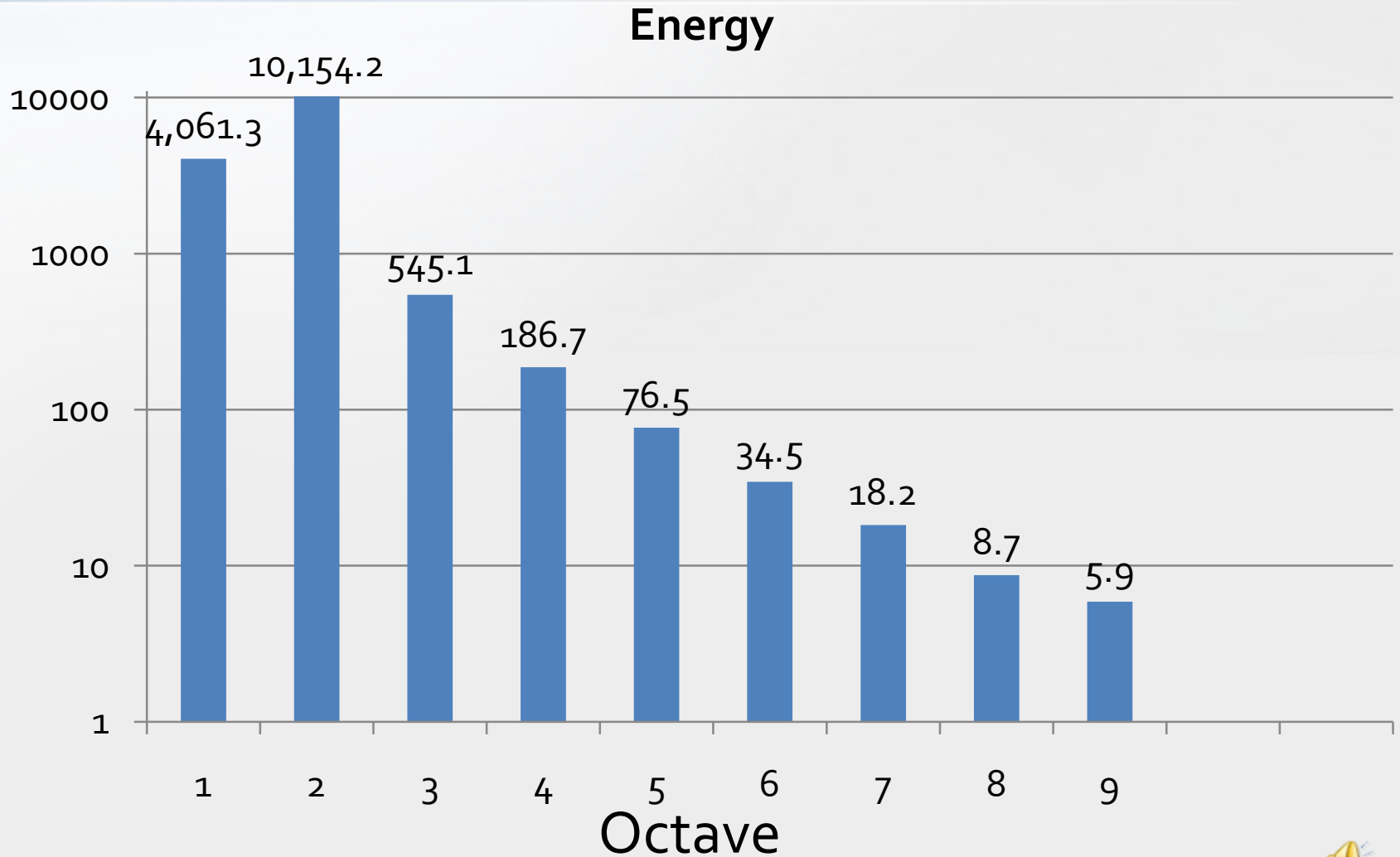
Stratified Cross-Validation

- Same as cross-validation, except that the folds are chosen so that they contain equal proportions of labels.

Spectral Bands



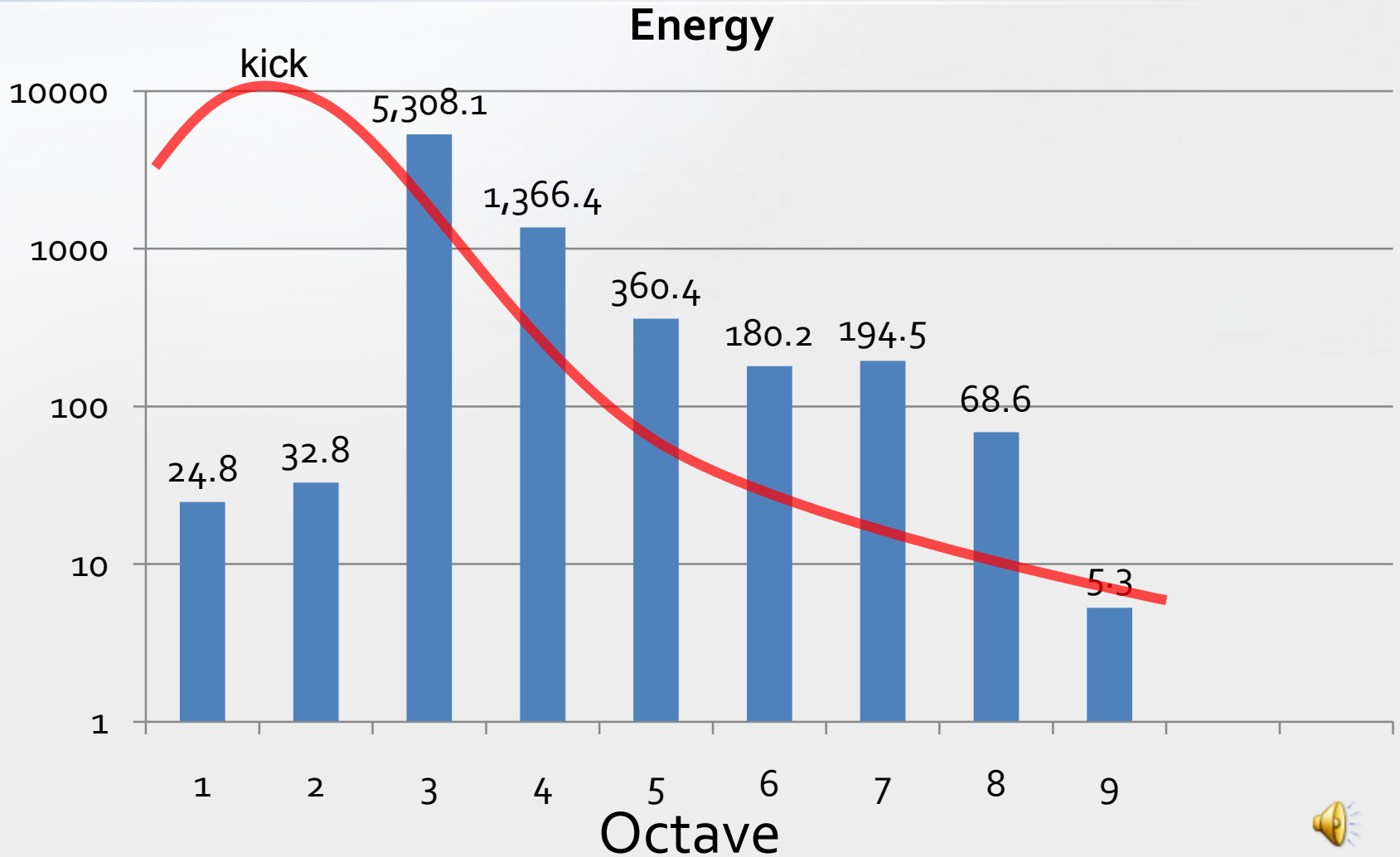
Frame 1



Features – Frame 1

Frame	ZC R	Centroid	BW	Skew	Kurtosis	E1	E2	E3	E4	E5	E6	E7	E8	E9
1	9	2.8kHz	5kHz	2.2	6.7	4000	10100	545	187	77	35	18	9	6

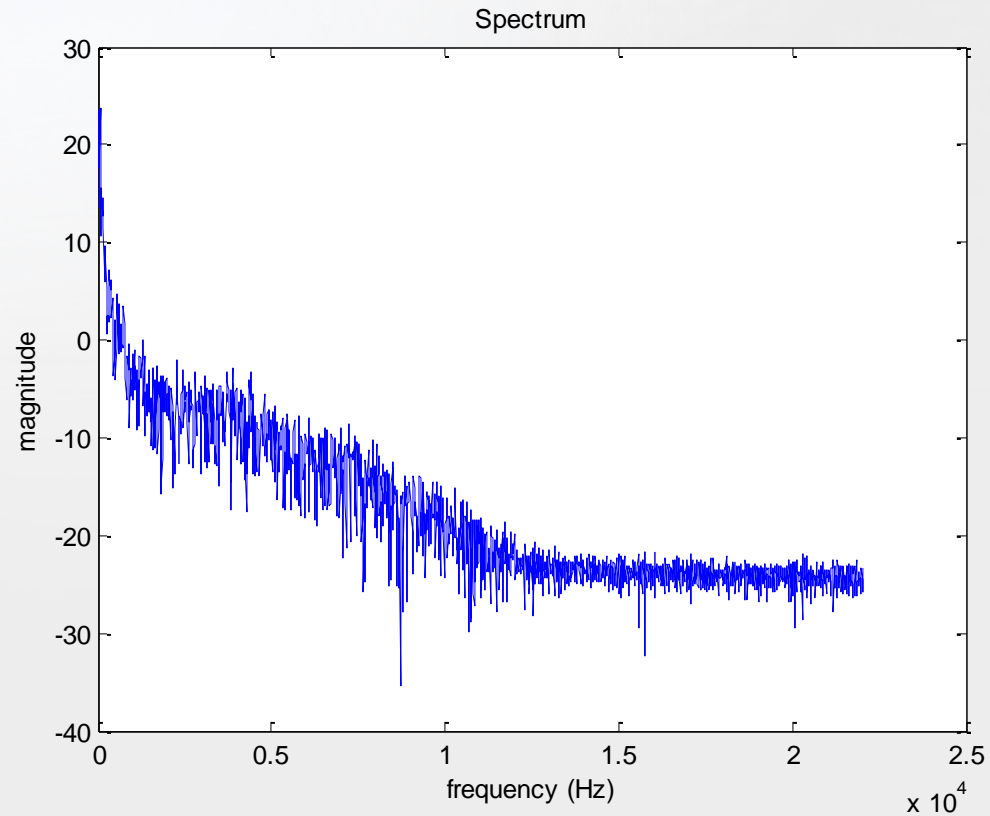
Frame 2

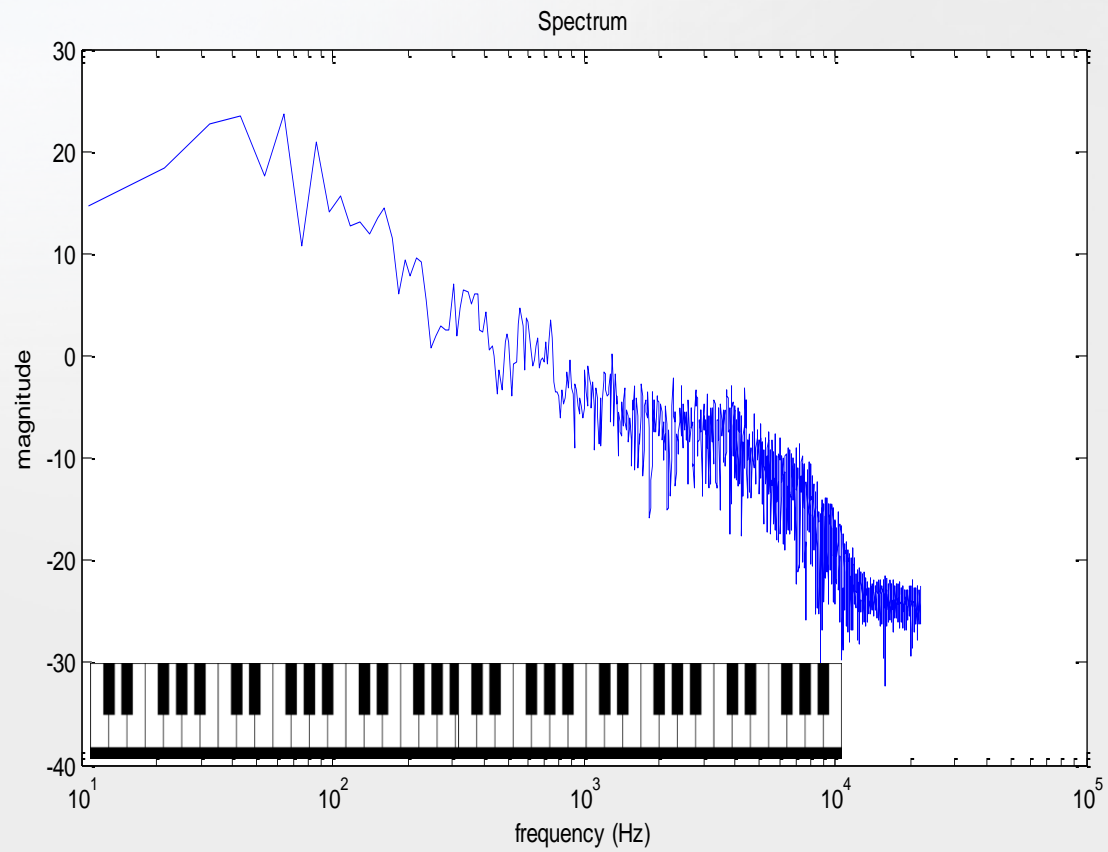


Features : SimpleLoop.wav

Frame	ZC R	Centroid	BW	Skew	Kurtosis	E1	E2	E3	E4	E5	E6	E7	E8
1	9	2.8kHz	5kHz	2.2	6.7	4000	10100	545	187	77	35	18	9
2	423	3.1kHz	4kHz	2	7.2	24	33	5300	1366	360	180	194	68

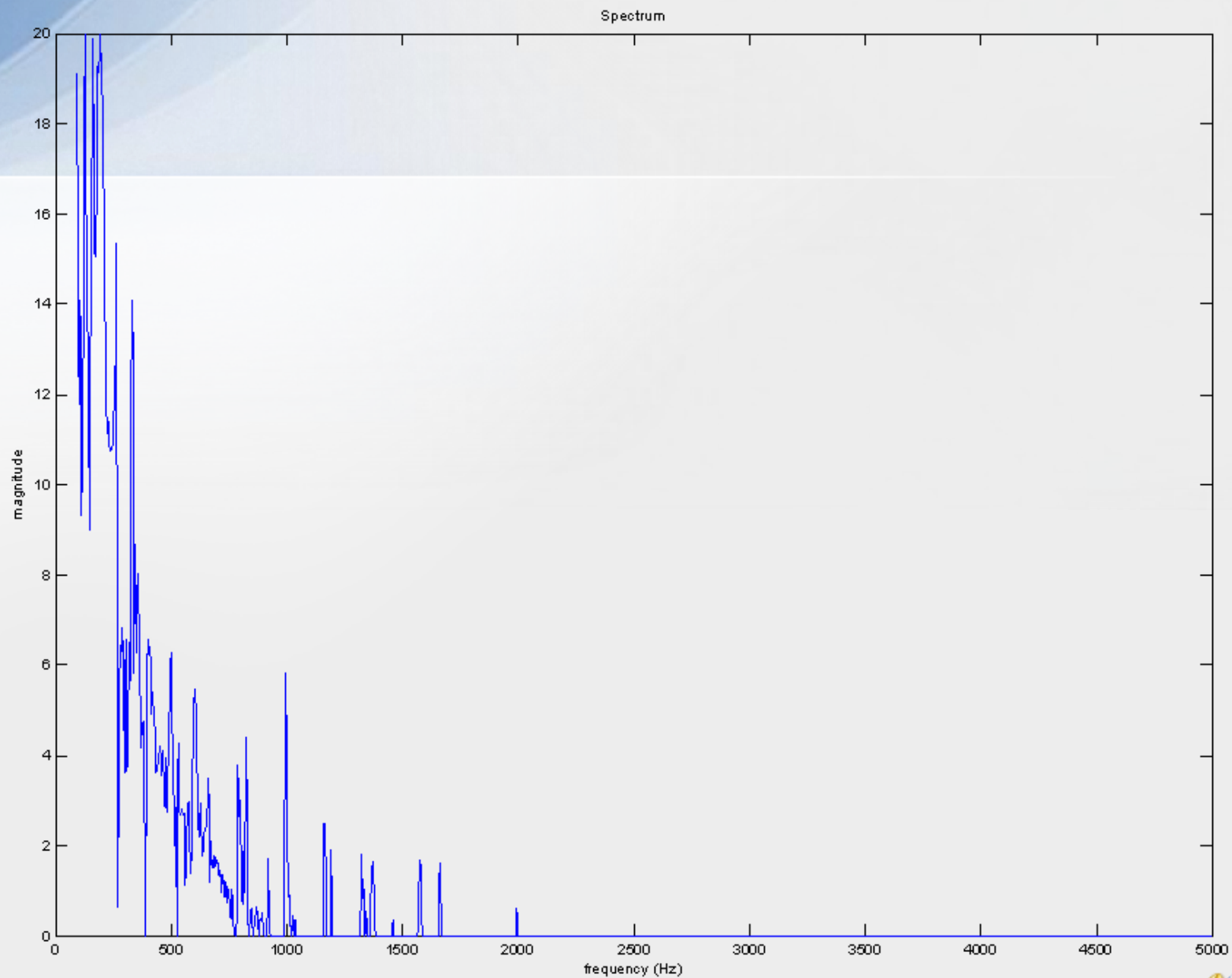
Log Spectrogram

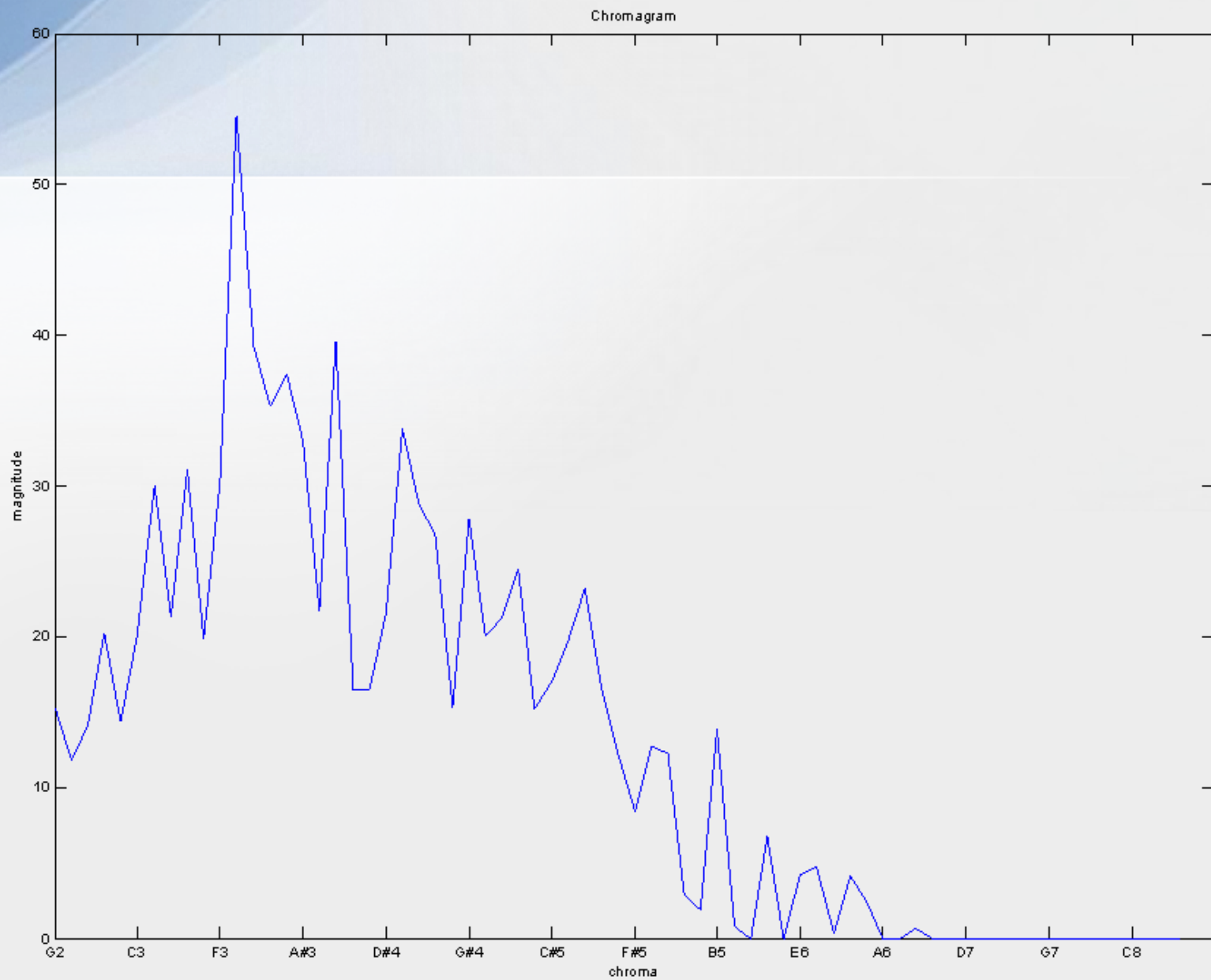


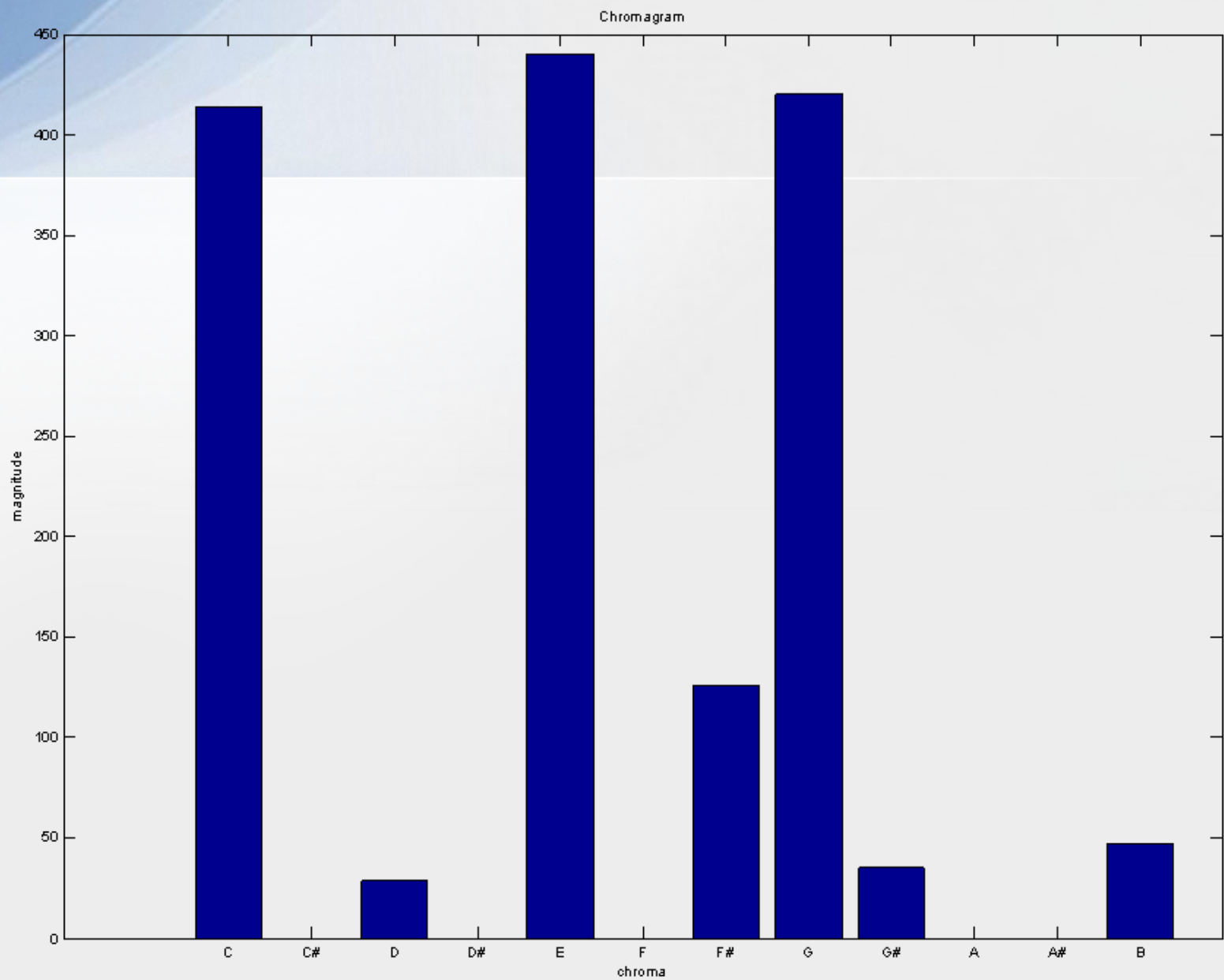


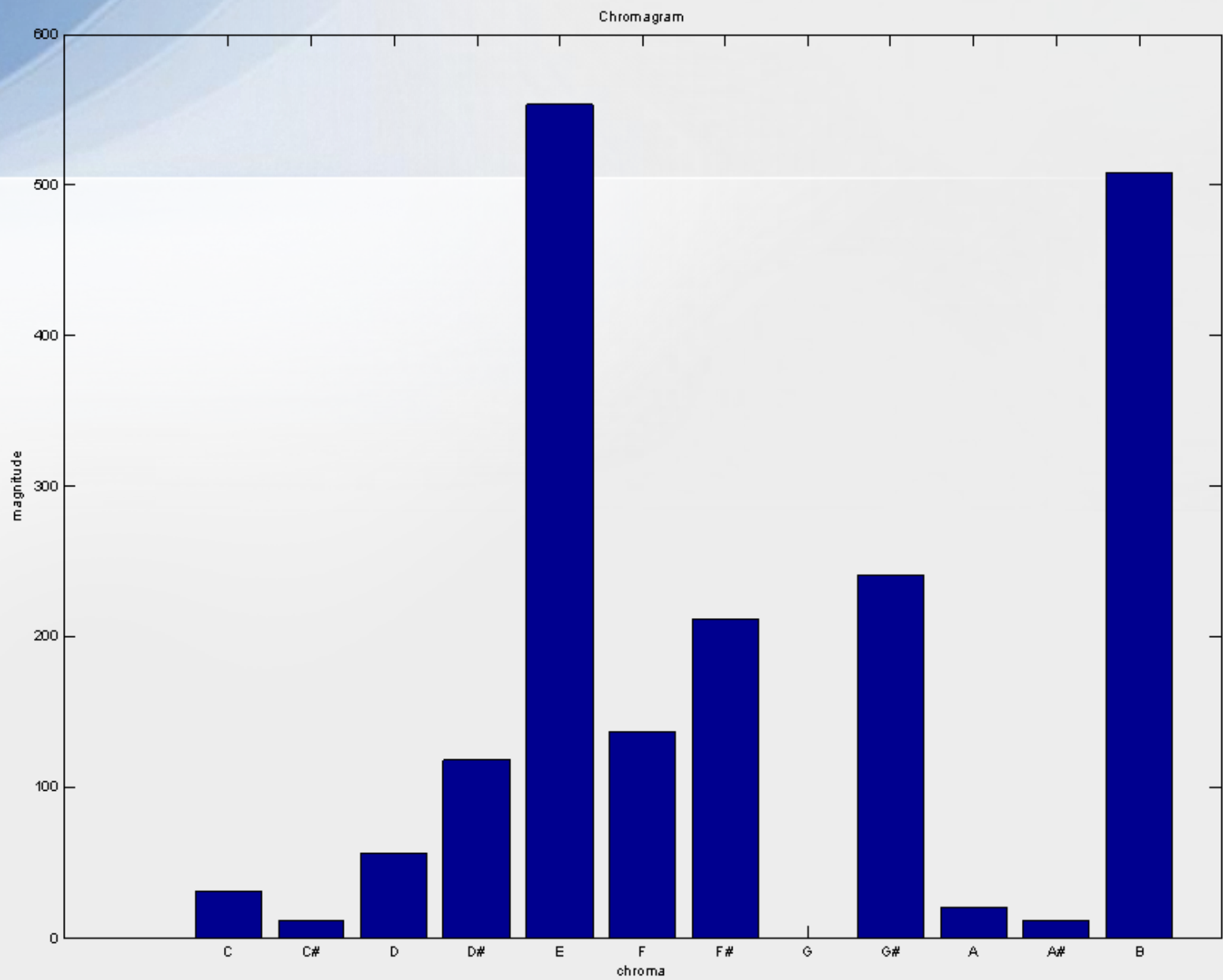
Chroma Bins



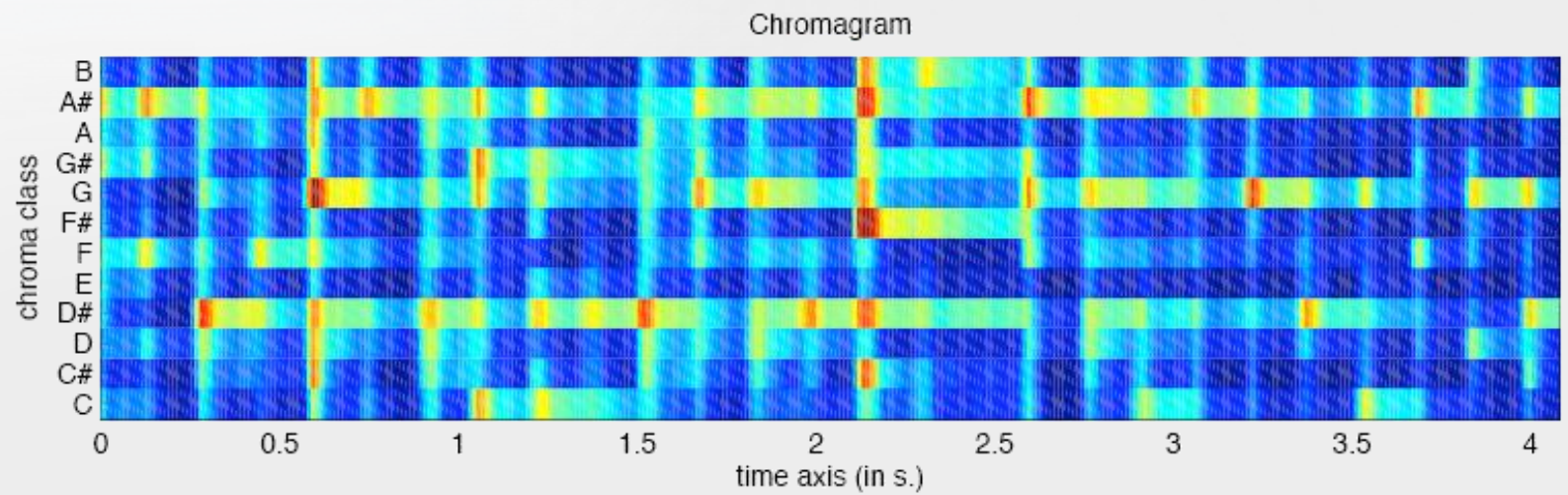




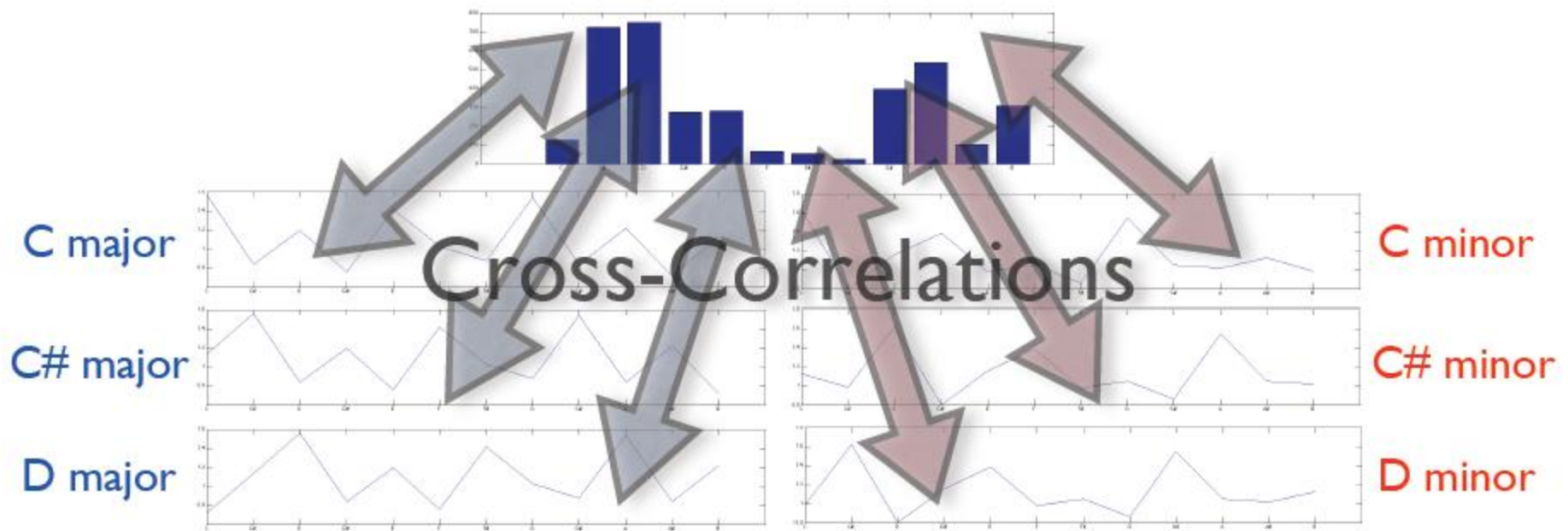




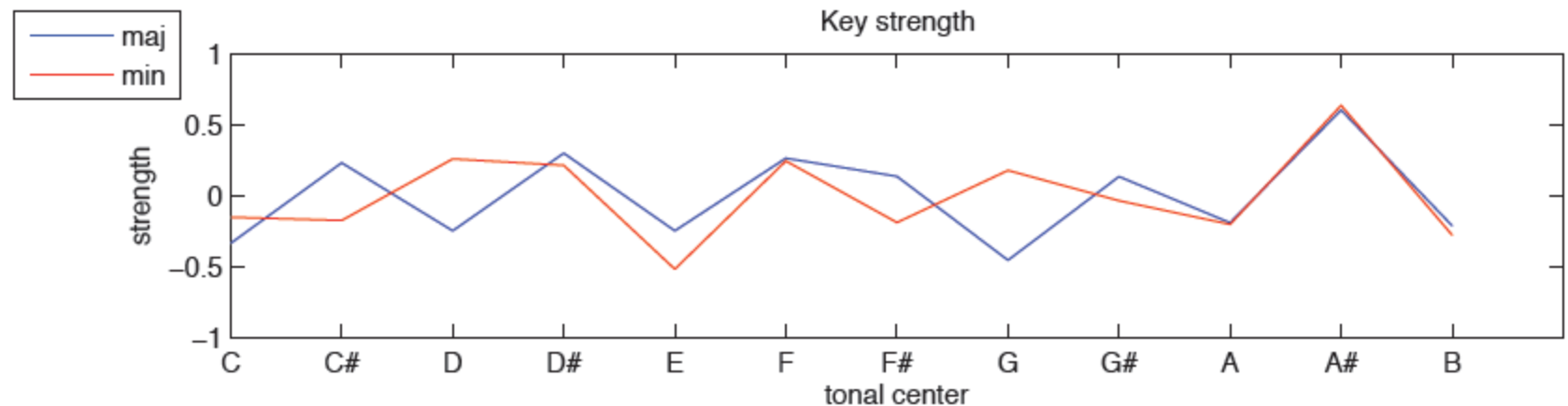
EXAMPLE



Picture courtesy: Olivier Lartillot



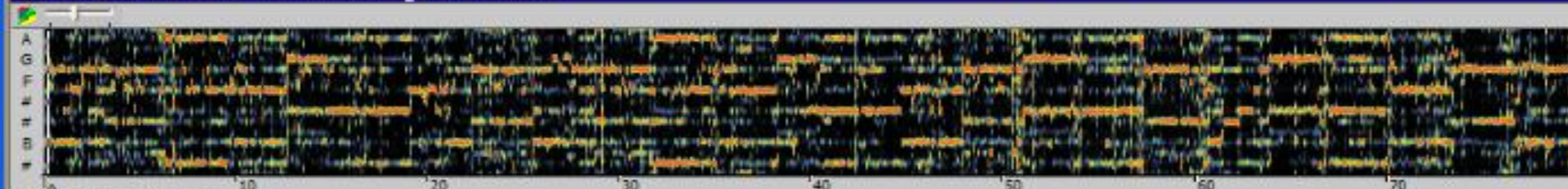
The resulting graph indicate the cross-correlation score for each different tonality candidate.



Input sound 1 - C:\Datos\Emilia\TestSounds\hotel_california-short.wav



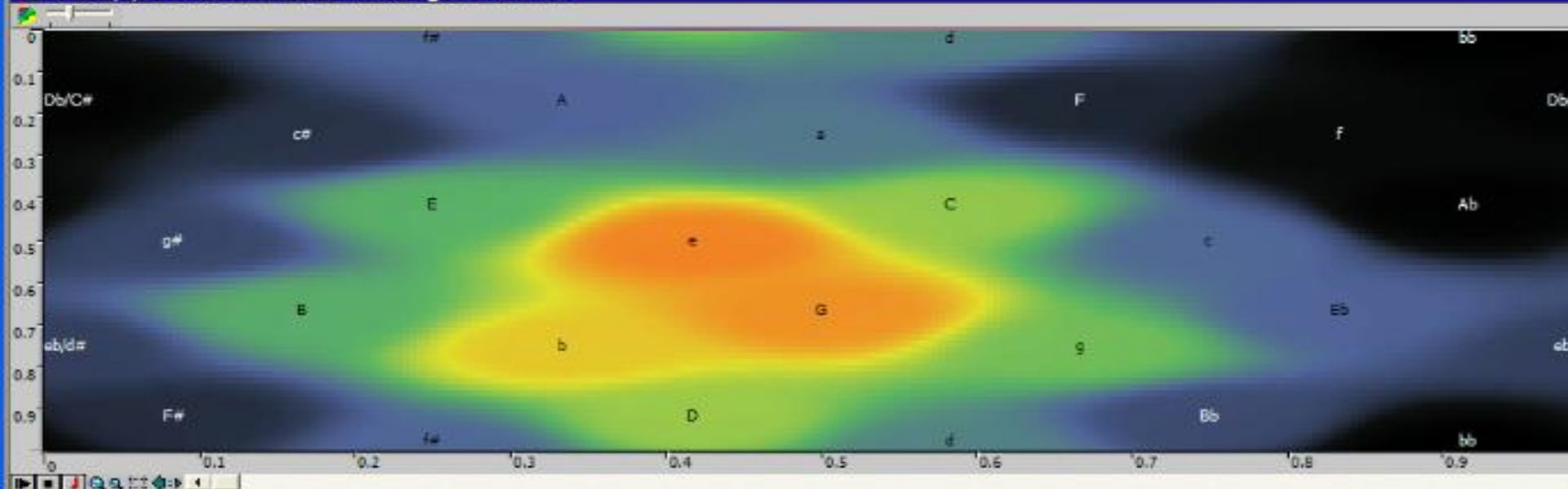
HPCPGrama 1 - C:\Datos\Emilia\TestSounds\hotel_california-short.wav

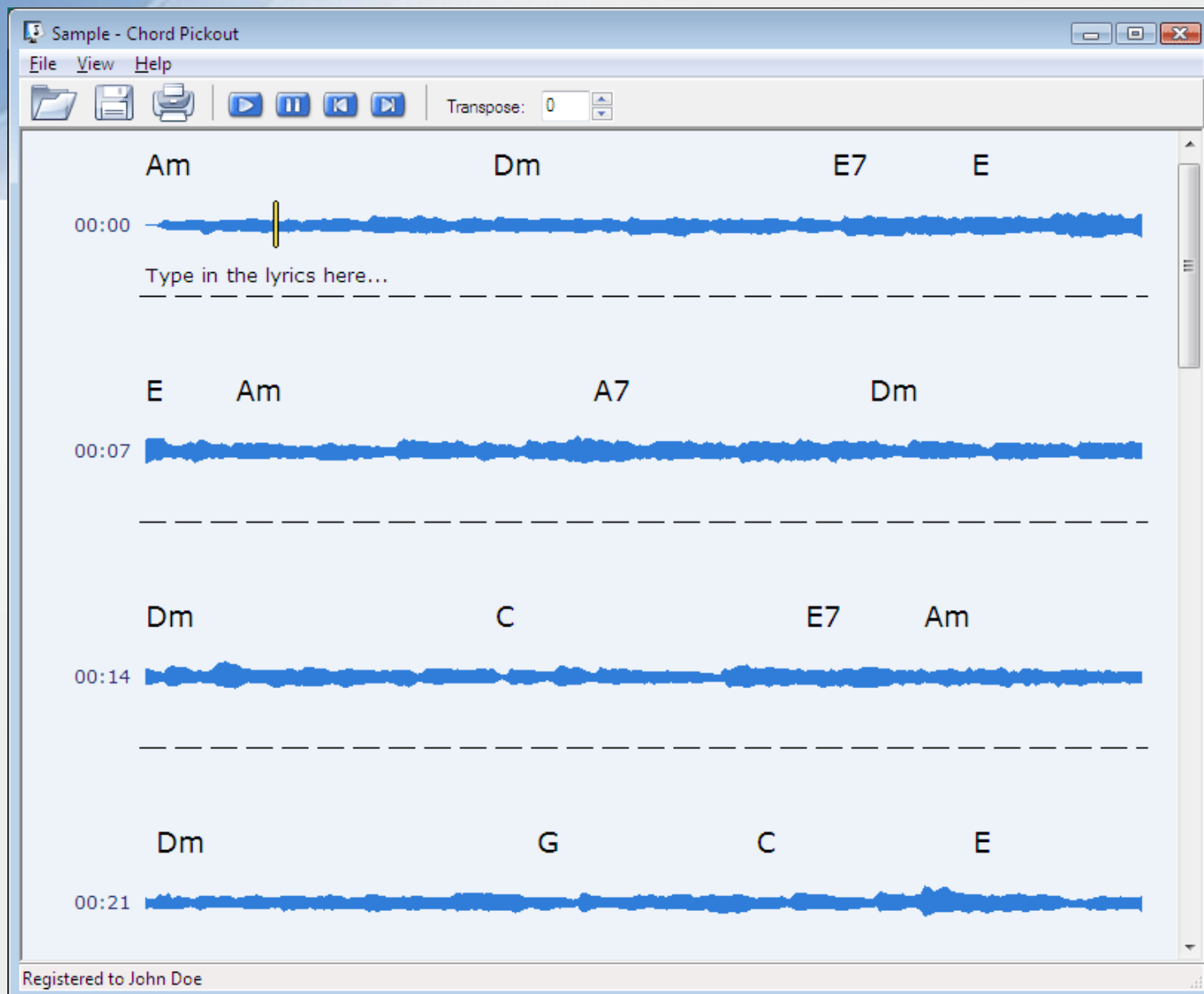


HPCPMean 1 - C:\Datos\Emilia\TestSounds\hotel_california-short.wav

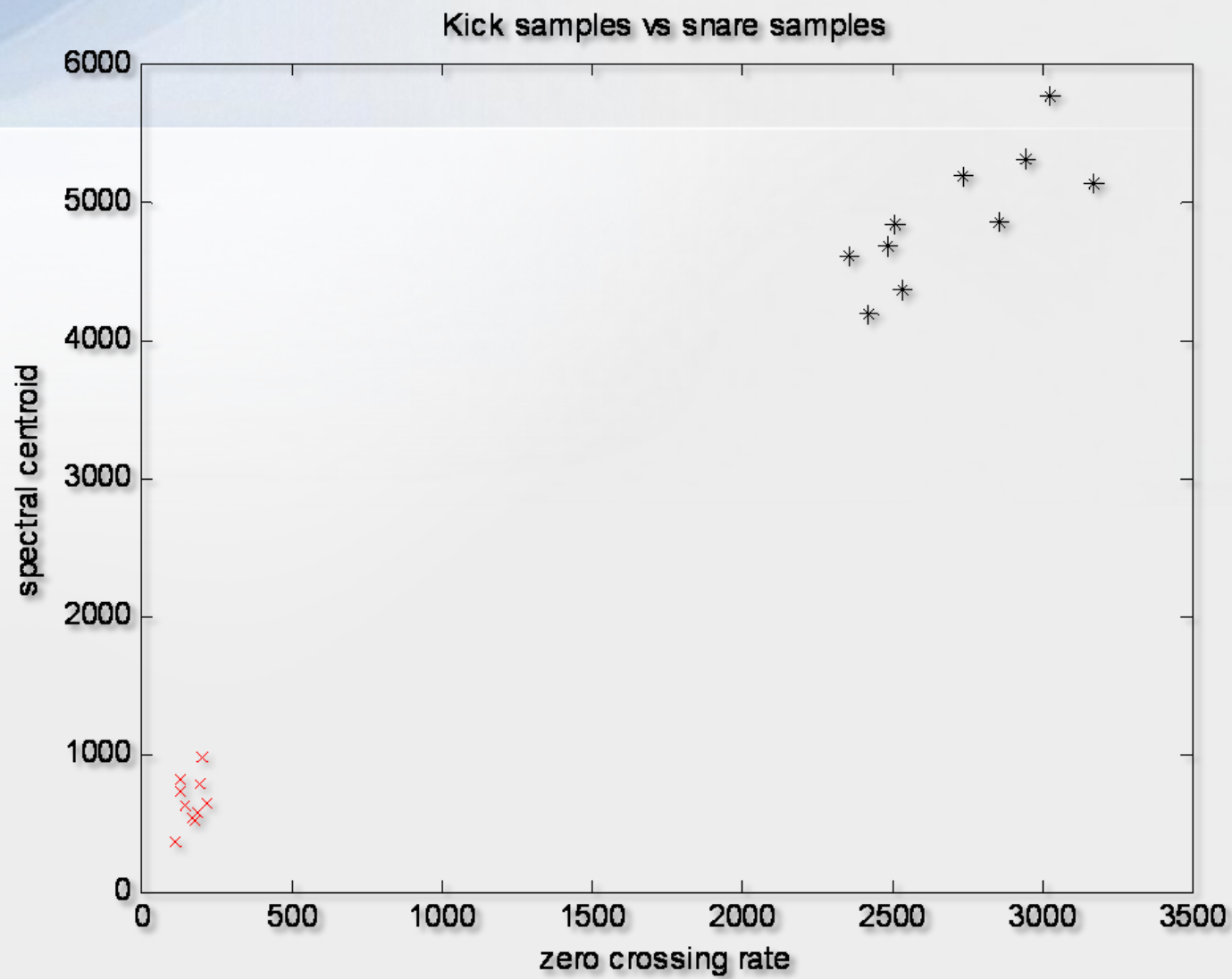


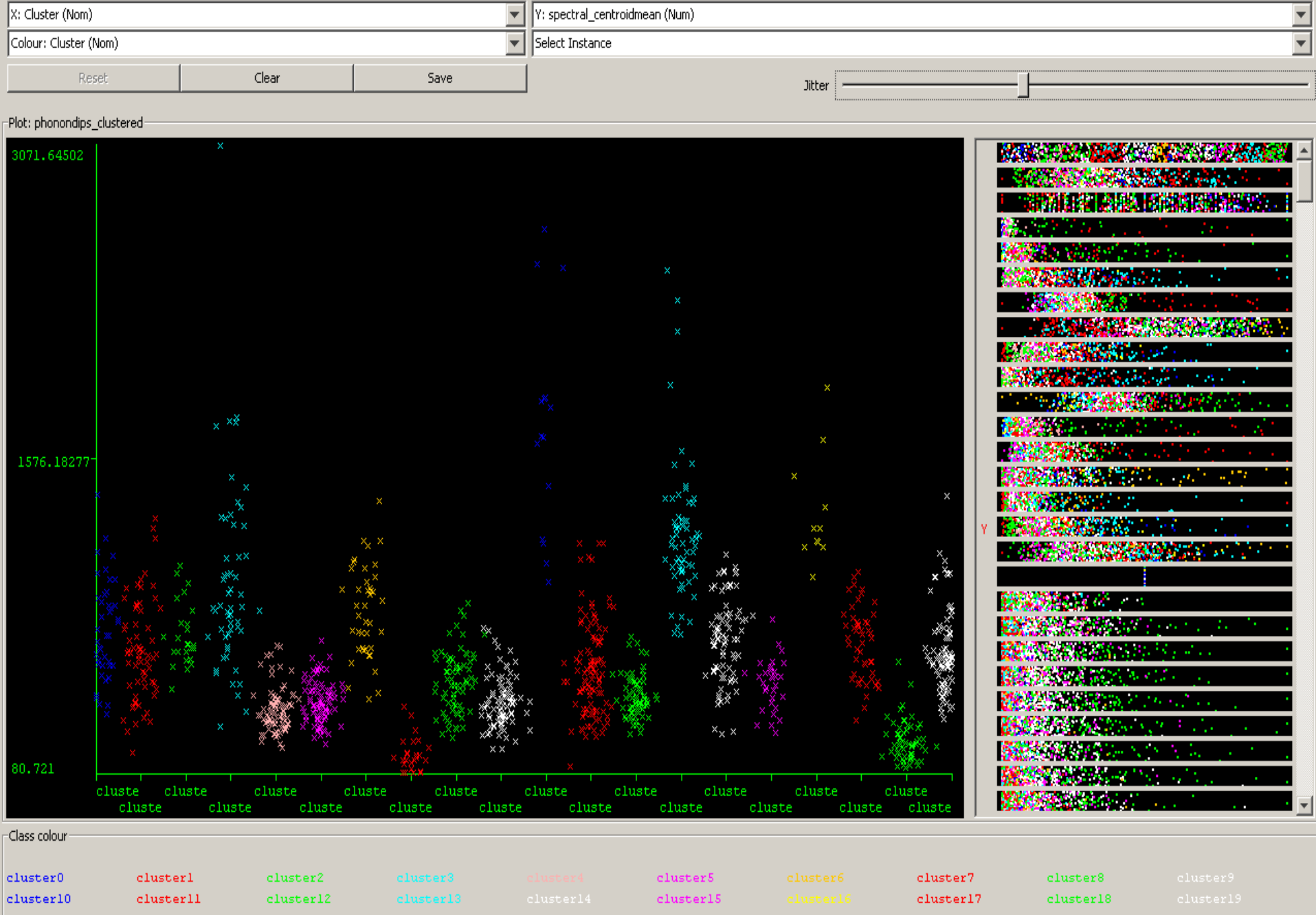
HPCPKeySpace 1 - C:\Datos\Emilia\TestSounds\hotel_california-short.wav





- <http://www.chordpickout.com/index.html>





Day 3 Lab

- Get your Lab 2 working
 - Make sure that training data = 100% accurate
 - Try the test snares and test kicks
 - Write down your accuracy and parameters
 - Change the number of features
 - Add or replace current features with different values
 - (e.g., mirbrightness, mirrolloff)
- Demo - tonality
- Demo – tempo