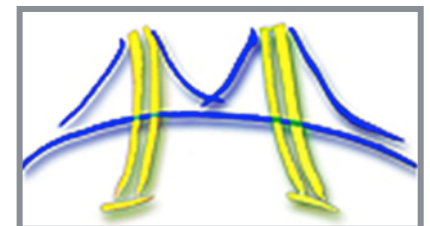# Toward Live Drum Separation Using Probabilistic Spectral Clustering Based on the Itakura-Saito Divergence

**Eric Battenberg**

**UC Berkeley, Dept. of EECS**

**Parallel Computing Laboratory**
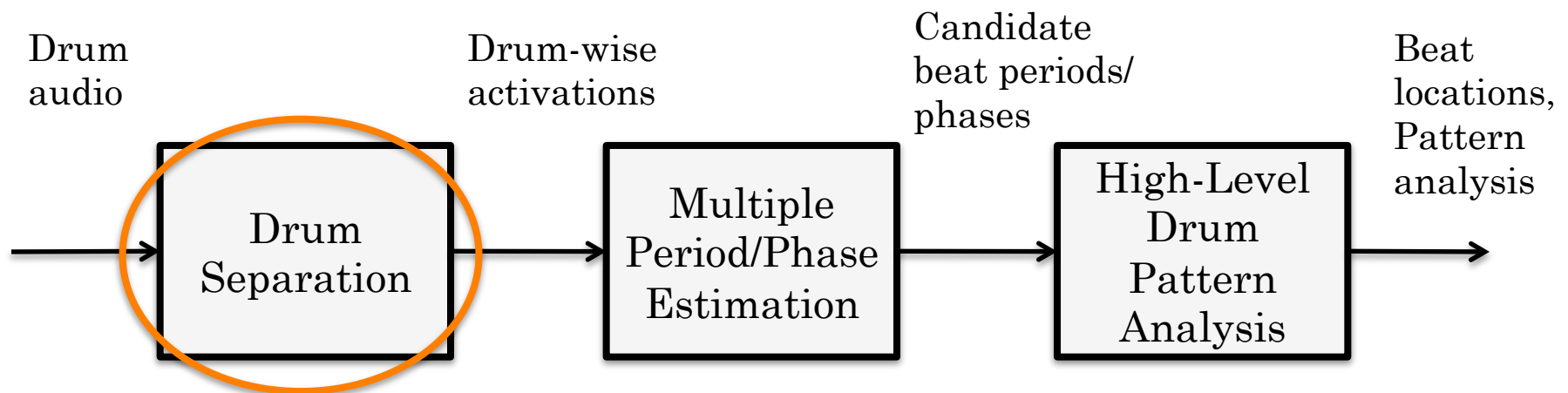
**CNMAT**

# TOWARD COMPREHENSIVE RHYTHMIC UNDERSTANDING

- Or "Live Drum Understanding"
- <u>Goal</u>: Go beyond simple beat tracking and provide context-aware, instrument-aware information in real-time, e.g.
  - "This rhythm is in 5/4 time"
  - "This drummer is playing syncopated notes on the hi-hat"
  - "The ride cymbal pattern has a swing feel"
  - "This is a Samba rhythm"

# LIVE DRUM UNDERSTANDING SYSTEM

Drum audio → **Drum Separation** → Drum-wise activations → **Multiple Period/Phase Estimation** → Candidate beat periods/ phases → **High-Level Drum Pattern Analysis** → Beat locations, Pattern analysis

- This work.
- Gamma Mixture Model training of drum templates
- Non-negative decomposition onto templates.

- HMM-based multi-hypothesis beat tracking.

- Statistical deep learning of drum patterns
- Stacked Conditional Restricted Boltzmann Machines

✓ ✓

# REQUIREMENTS FOR DRUM SEPARATION

- Real-Time/Live operation
- Useful with <u>any</u> percussion setup.
  - Before a performance, we can quickly train the system for a particular percussion setup.
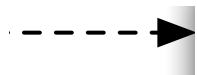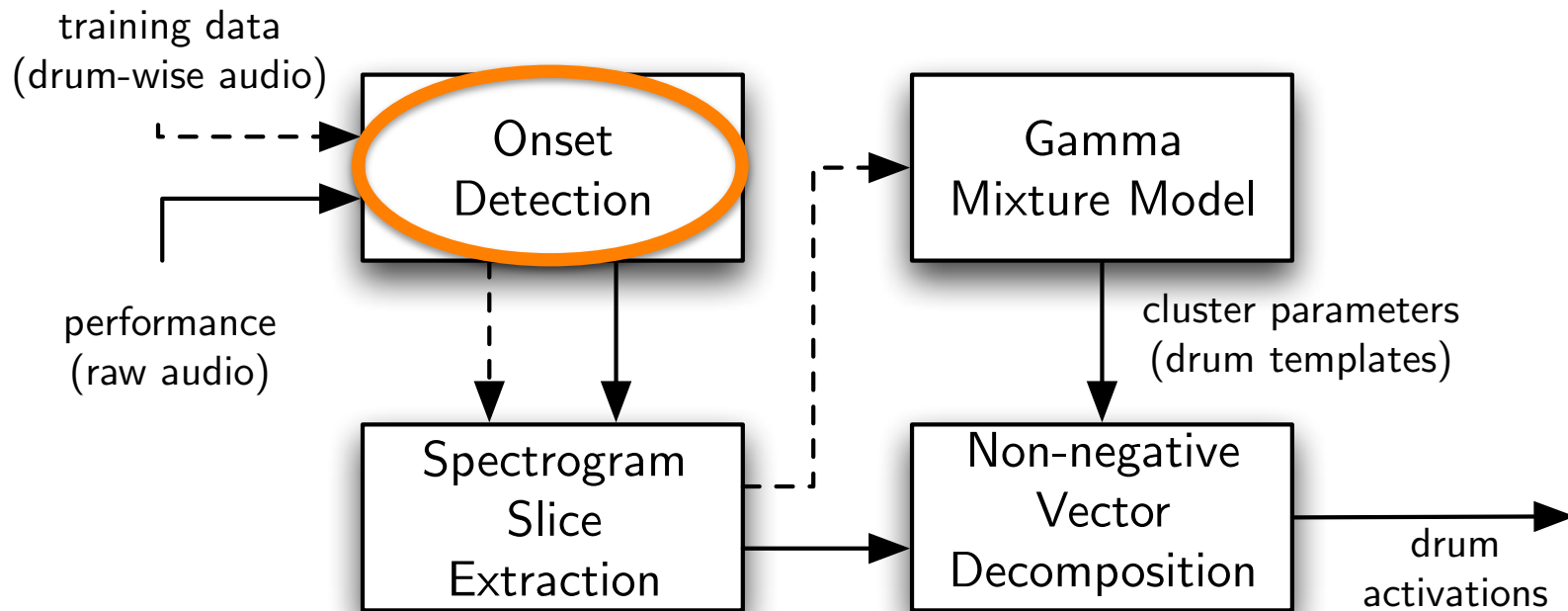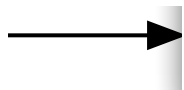
# THE PRIMARY TAKEAWAY

- Gamma Mixture Model
  - *For learning spectral drum templates.*
  - **Cheaper to train** than GMM
  - **More stable** than GMM
- Non-negative Vector Decomposition (NVD)
  - *For computing template activations from drum onsets.*
  - Learning **multiple templates per drum** improves separation.
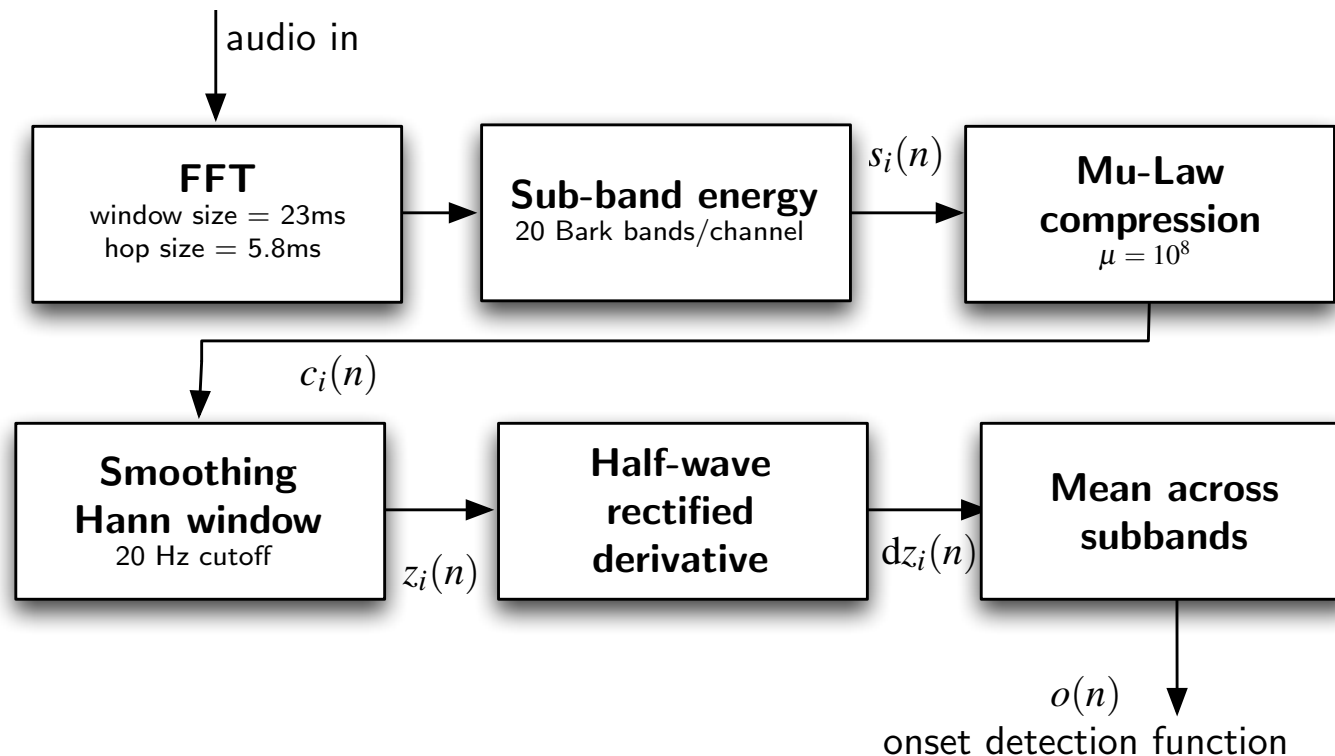  - The use of **"tail" templates** reduces false positives.
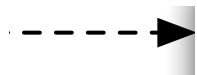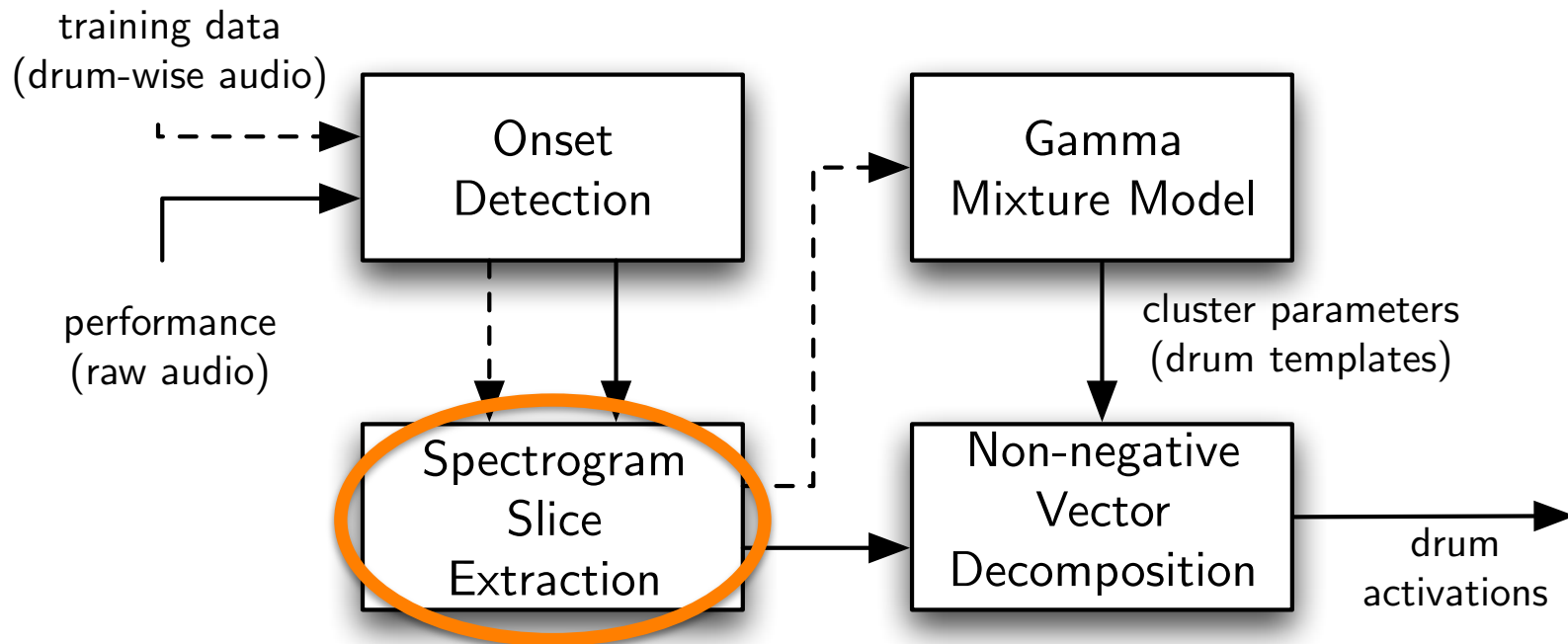
# DRUM SEPARATION SYSTEM

training data
(drum-wise audio)

performance
(raw audio)

**Onset Detection**

**Gamma Mixture Model**

cluster parameters
(drum templates)

**Spectrogram Slice Extraction**

**Non-negative Vector Decomposition**

drum activations
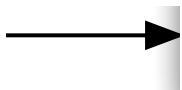
# ONSET DETECTION

- Detection function: Differentiated log-energy of multiple perceptual sub-bands.
- On 2400 drum strikes, our **adaptive threshold** achieves:
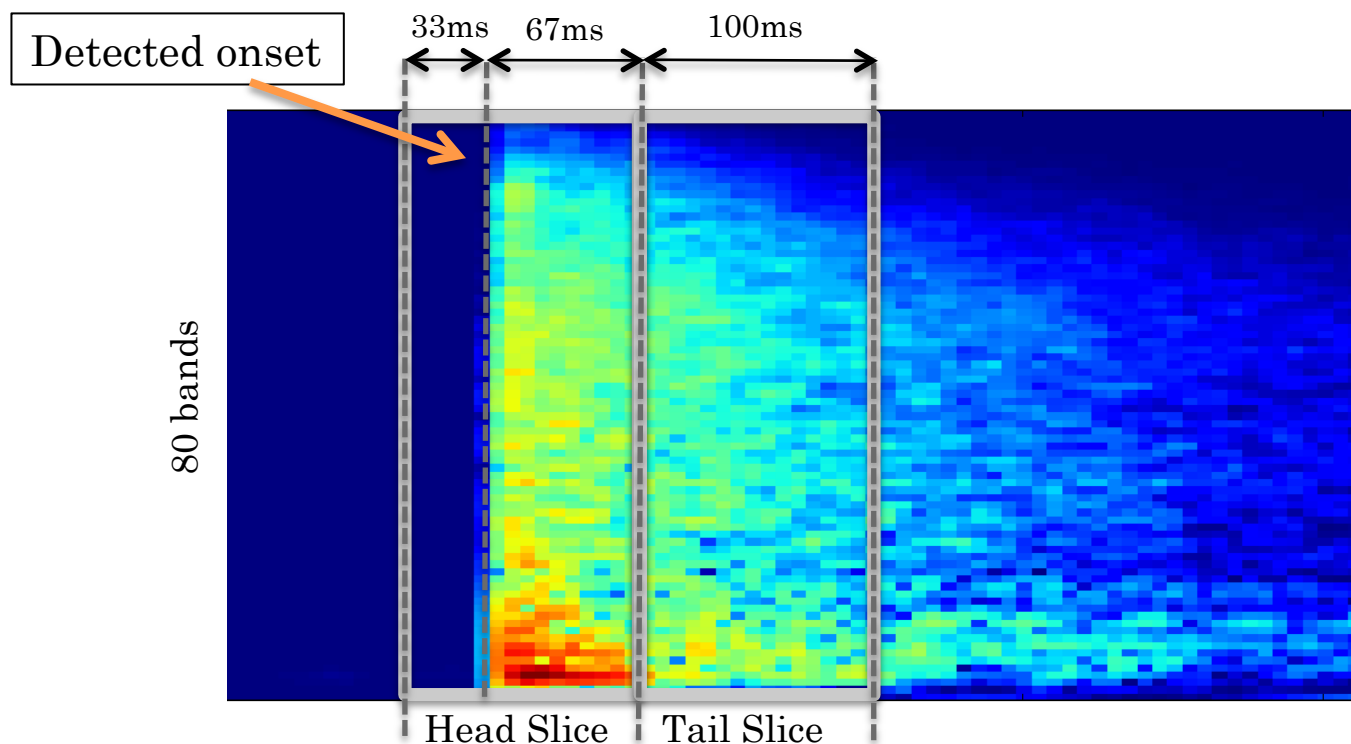  - **85% recall, 99.9% precision**.

audio in

| FFT | Sub-band energy | Mu-Law |
|---|---|---|
| window size = 23ms<br>hop size = 5.8ms | 20 Bark bands/channel | compression<br>$\mu = 10^8$ |

$s_i(n)$

$c_i(n)$

| Smoothing<br>Hann window | Half-wave<br>rectified | Mean across |
|---|---|---|
| 20 Hz cutoff | derivative | subbands |

$z_i(n)$

$\mathrm{d}z_i(n)$

$o(n)$
onset detection function

# DRUM SEPARATION SYSTEM

training data
(drum-wise audio)

performance
(raw audio)

Onset
Detection

Gamma
Mixture Model

cluster parameters
(drum templates)

Spectrogram
Slice
Extraction

Non-negative
Vector
Decomposition

drum
activations
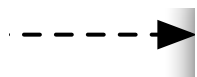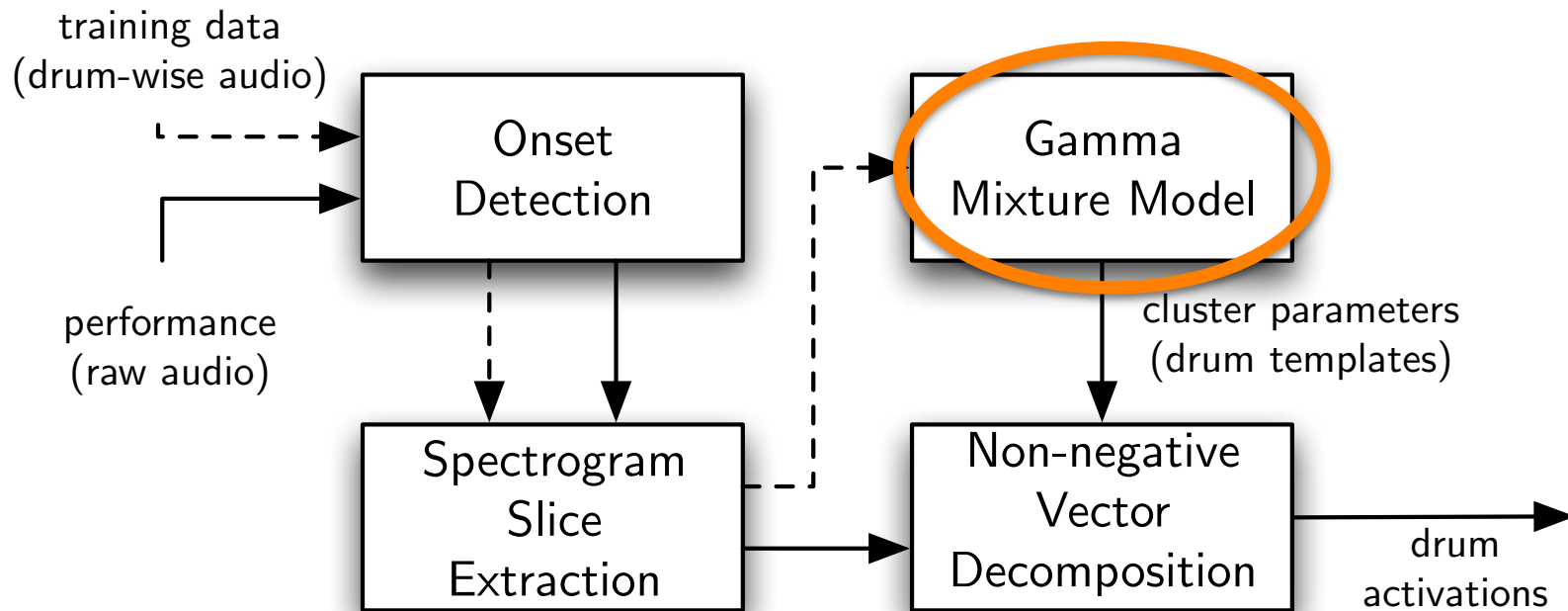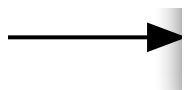
# SPECTROGRAM SLICES

- Extracted at onsets.
- Each slice contains 100ms (~17 frames) of audio
- 80 bark-spaced bands per channel [Battenberg 2008]
- During training, both "head" and "tail" slices are extracted.
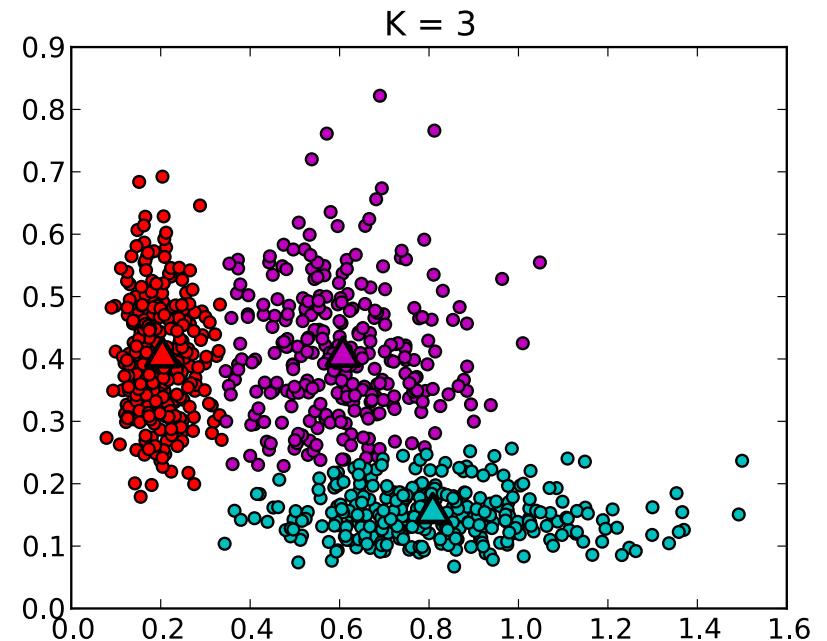  - Tail templates serve as decoys during non-negative vector decomposition.



Detected onset

33ms    67ms    100ms

80 bands

Head Slice    Tail Slice

# DRUM SEPARATION SYSTEM



training data
(drum-wise audio)

Onset
Detection

Gamma
Mixture Model

performance
(raw audio)

cluster parameters
(drum templates)

Spectrogram
Slice
Extraction

Non-negative
Vector
Decomposition

drum
activations

# TRAINING DRUM TEMPLATES

- Instead of taking an "average" of all training slices for a single drum…

- Cluster them and use the cluster centers as the drum templates.
  - This gives us multiple templates per drum…
  - Which helps represent the variety of sounds that can be made by a single drum.



K = 3

# CLUSTERING USING MIXTURE MODELS

- Train using the Expectation-Maximization (EM) algorithm.

- Gaussian Mixture Model (GMM)
  - Requires expensive/unstable covariance matrices
  - Enforces a <u>Euclidean distance</u> measure.

$$d_{Euc}(X,Y) = \int_{\omega} (X(\omega) - Y(\omega))^2 \, d\omega$$

- Gamma Mixture Model
  - Single mean vector per component
  - Enforces an <u>Itakura-Saito (IS) distance</u> measure
    - A scale-invariant perceptual distance between audio spectra.

$$d_{IS}(X,Y) = \int_{\omega} \left[ \frac{X(\omega)}{Y(\omega)} - \log \frac{X(\omega)}{Y(\omega)} - 1 \right] d\omega$$
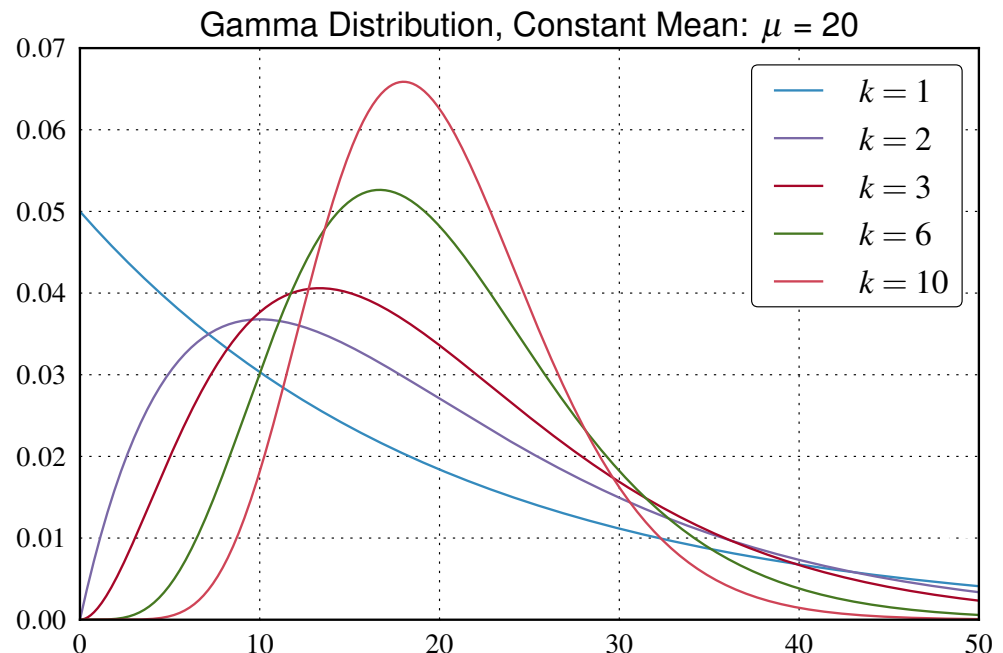
# GAMMA DISTRIBUTION

- Our mixture model is composed of gamma distributions.
- The gamma distribution models the sum of k independent exponential distributions.

$$p(y|\lambda,k) = y^{k-1}\frac{\lambda^k e^{-\lambda y}}{\Gamma(k)},$$
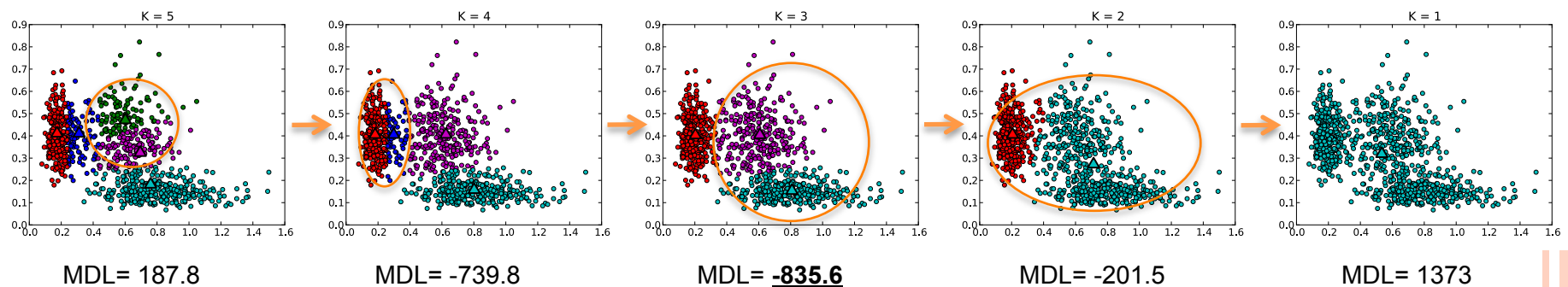
$$E[y] = \mu = k/\lambda$$

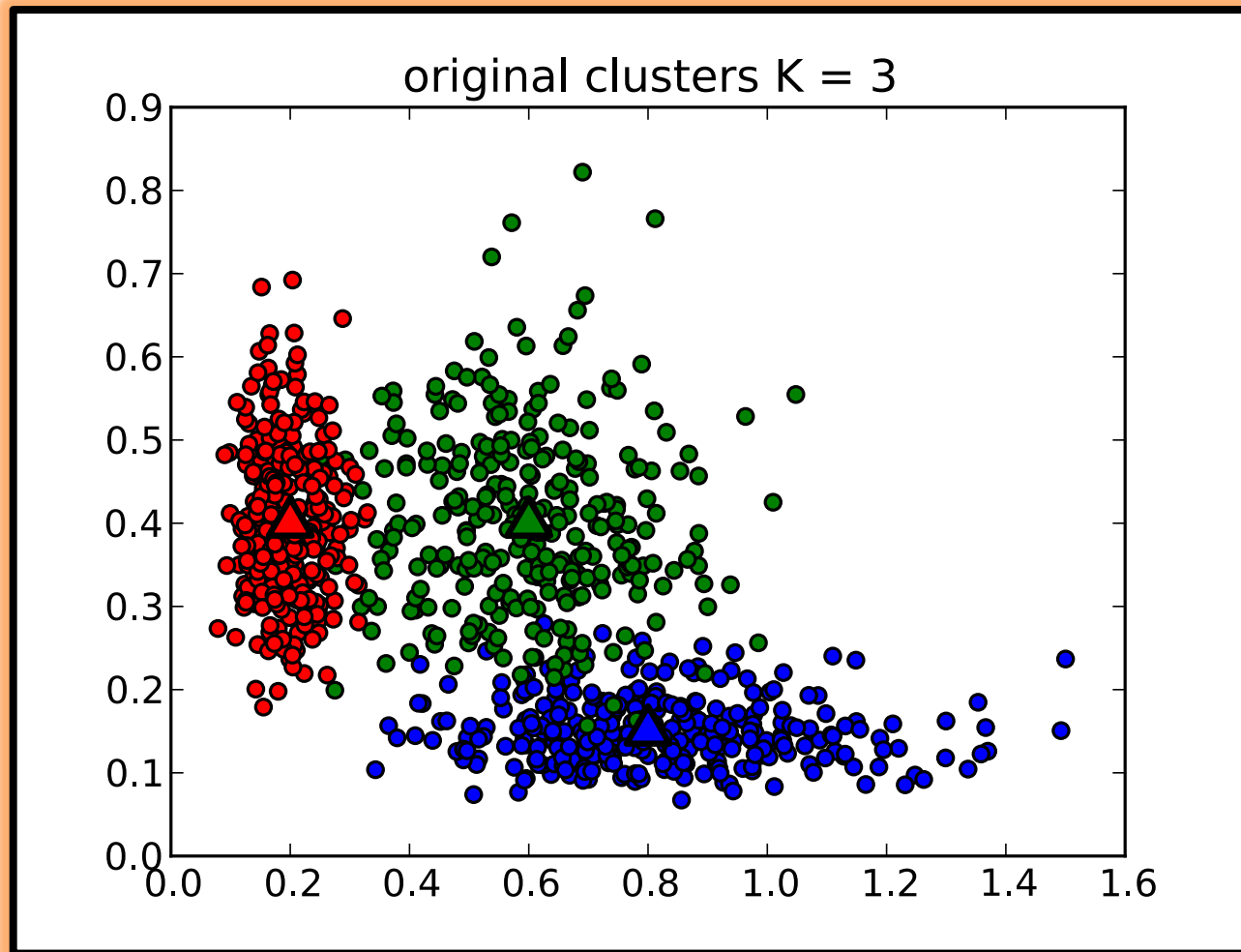$$\text{Var}[y] = \mu^2/k = k/\lambda^2$$

$$y \geq 0; \ \lambda, k > 0$$

Gamma Distribution, Constant Mean: $\mu = 20$

# AGGLOMERATIVE CLUSTERING

- *How many clusters to train?*
- We use Minimum Description Length (MDL) to choose the number of clusters.
  - Negative log-likelihood
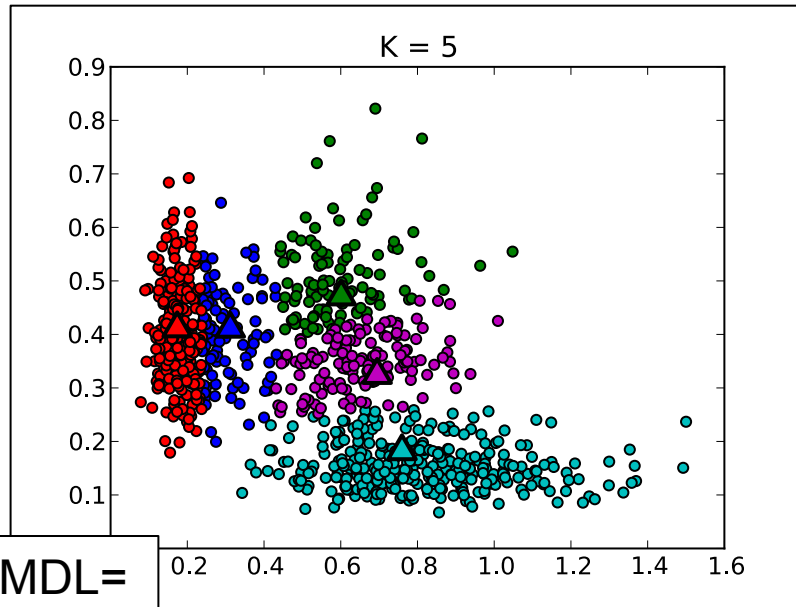  - + penalty term for number of clusters.



| MDL= 187.8 | MDL= -739.8 | MDL= **-835.6** | MDL= -201.5 | MDL= 1373 |

- 1. <u>Run EM</u> to convergence.
- 2. <u>Merge</u> the two most similar clusters.
- 3. <u>Repeat</u> 1,2 until we have a single cluster.
- 4. *Choose parameter set with smallest MDL.*

14

# AGGLOMERATIVE CLUSTERING WITH MDL
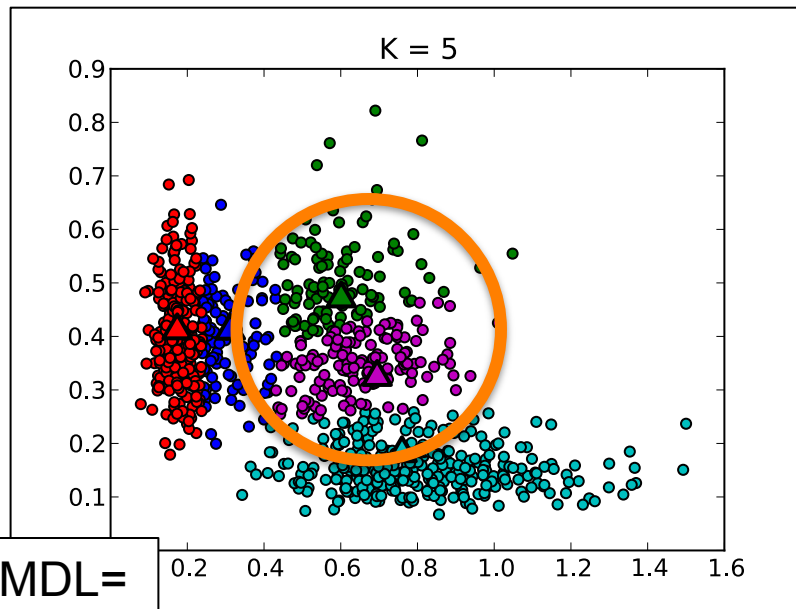


Ground Truth Mixture Data

# AGGLOMERATIVE CLUSTERING WITH MDL
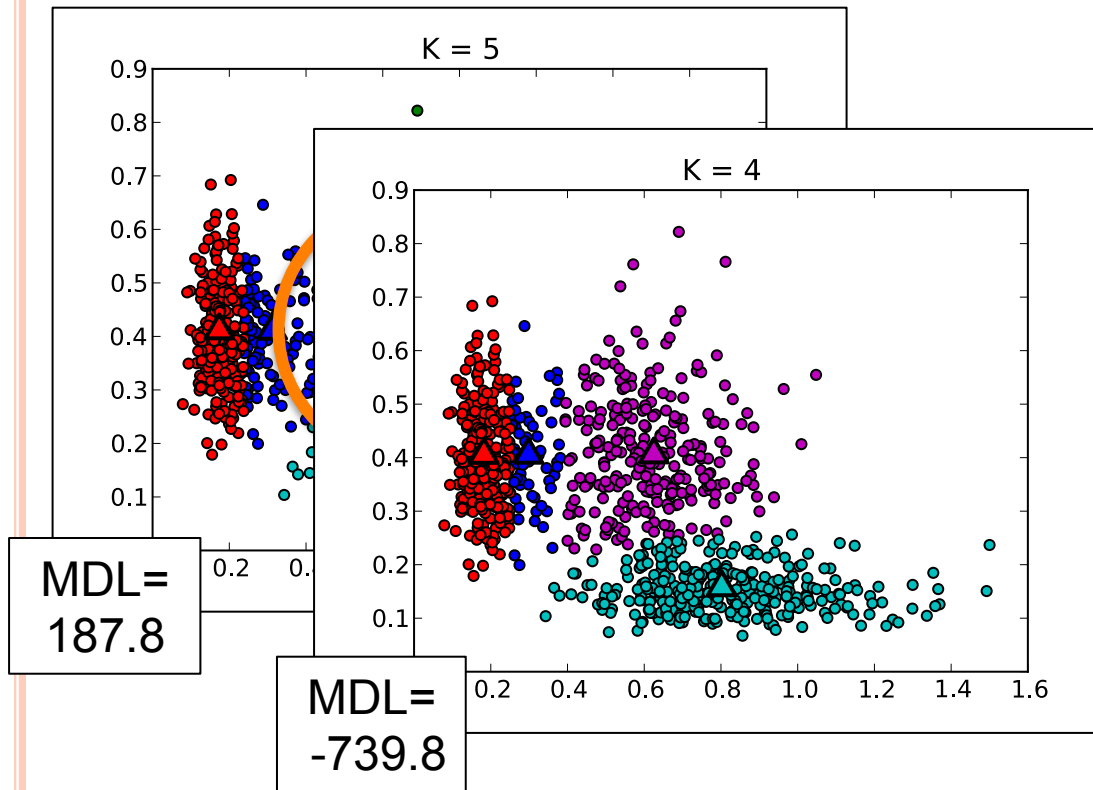


MDL=
187.8

# AGGLOMERATIVE CLUSTERING WITH MDL
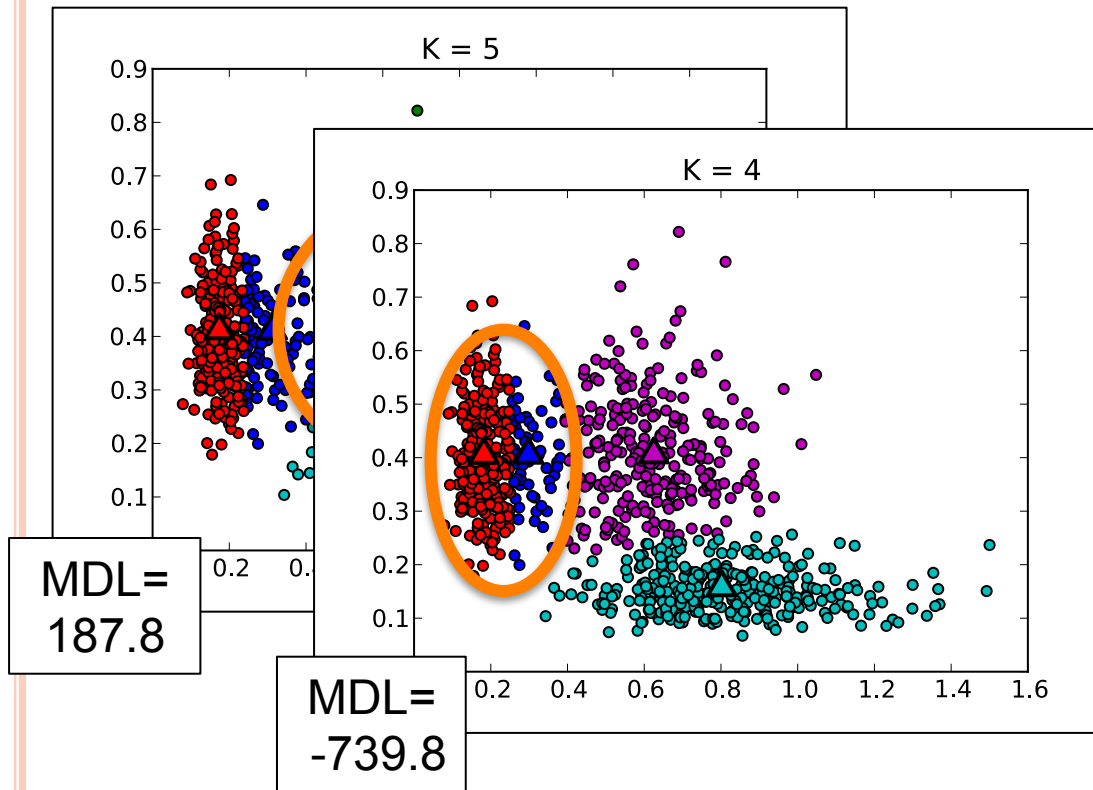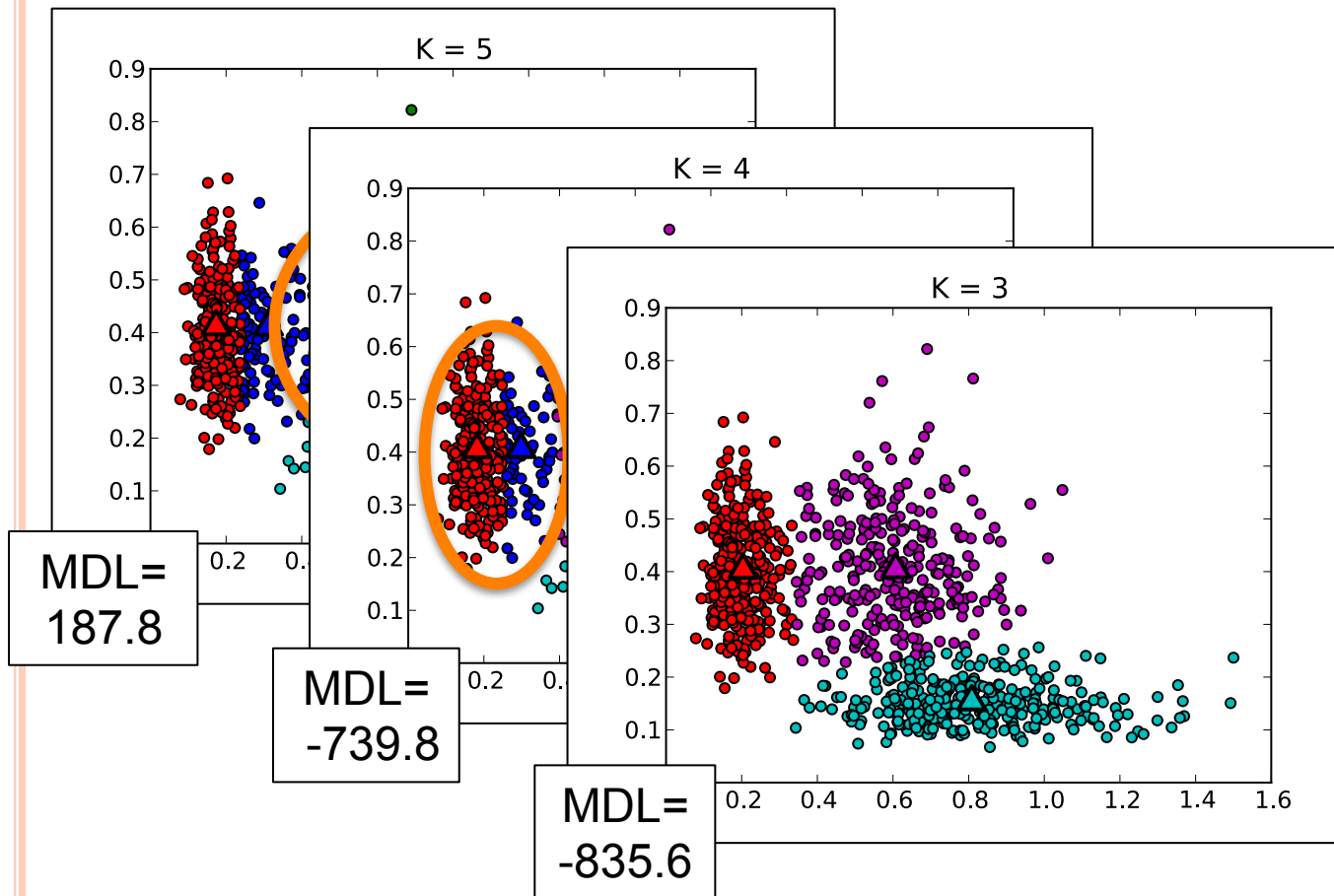


MDL= 187.8
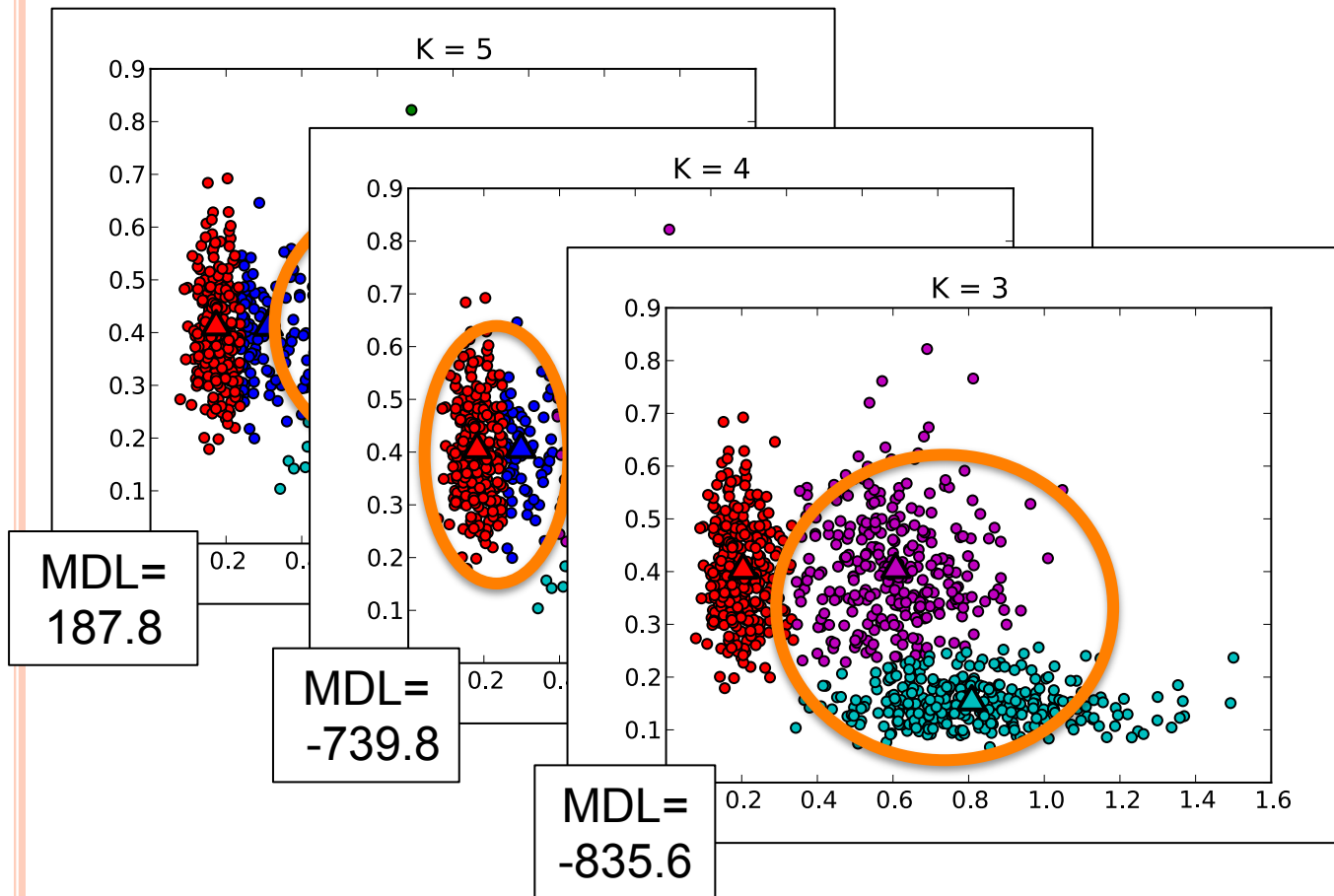
# Agglomerative Clustering with MDL

# AGGLOMERATIVE CLUSTERING WITH MDL

# AGGLOMERATIVE CLUSTERING WITH MDL



K = 5

MDL= 187.8

K = 4

MDL= -739.8

K = 3

MDL= -835.6

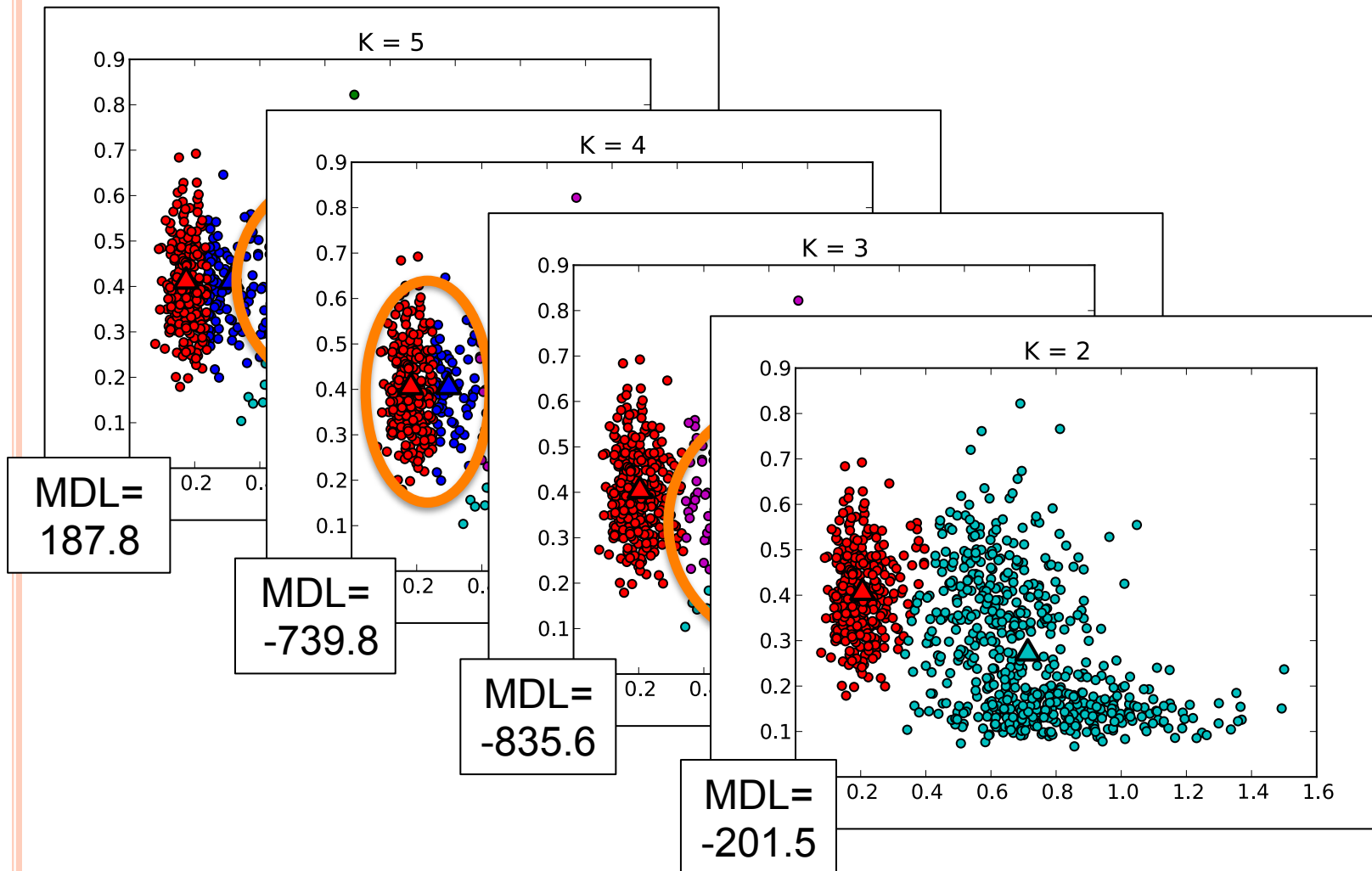# AGGLOMERATIVE CLUSTERING WITH MDL

# AGGLOMERATIVE CLUSTERING WITH MDL

# AGGLOMERATIVE CLUSTERING WITH MDL

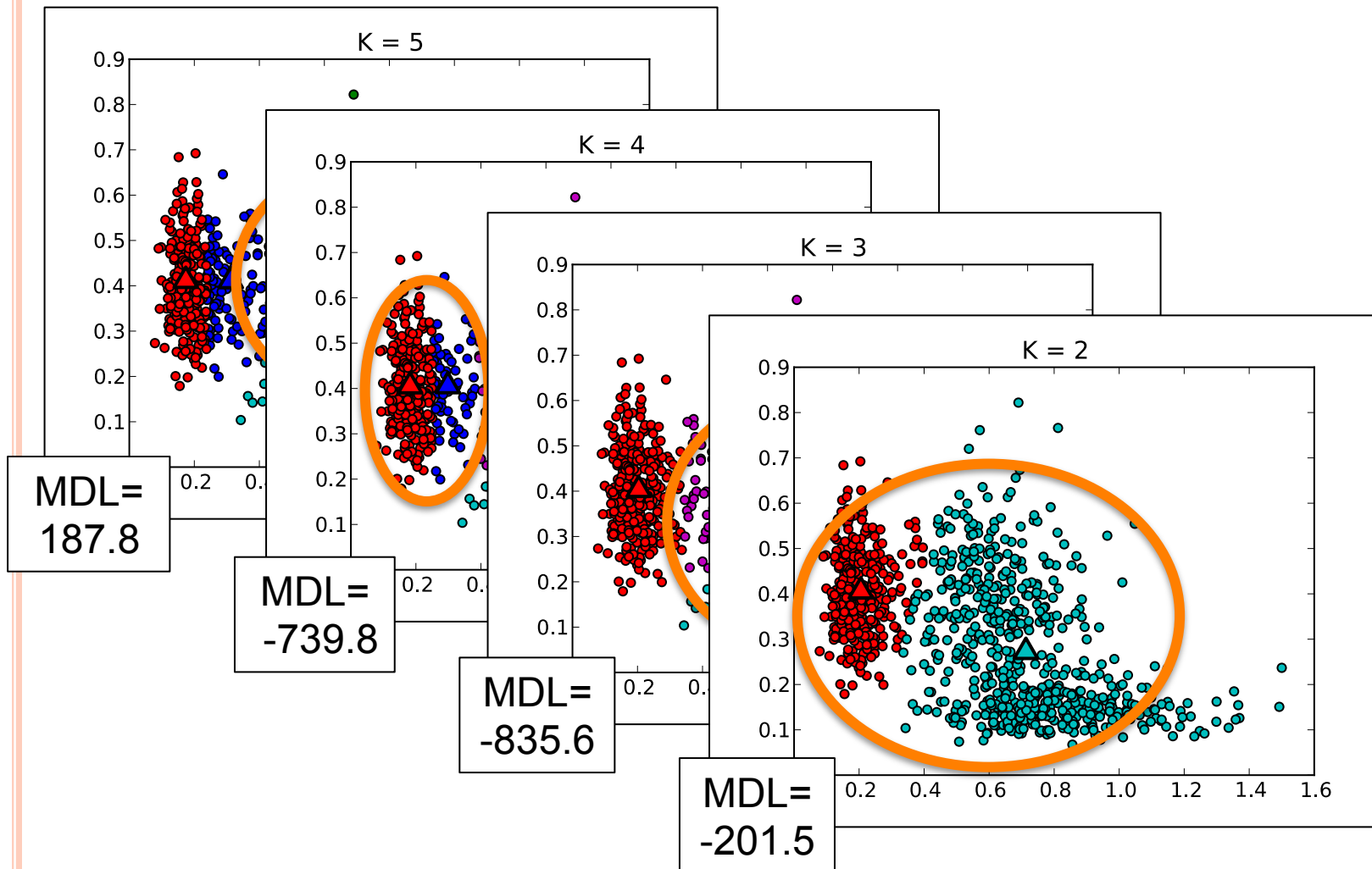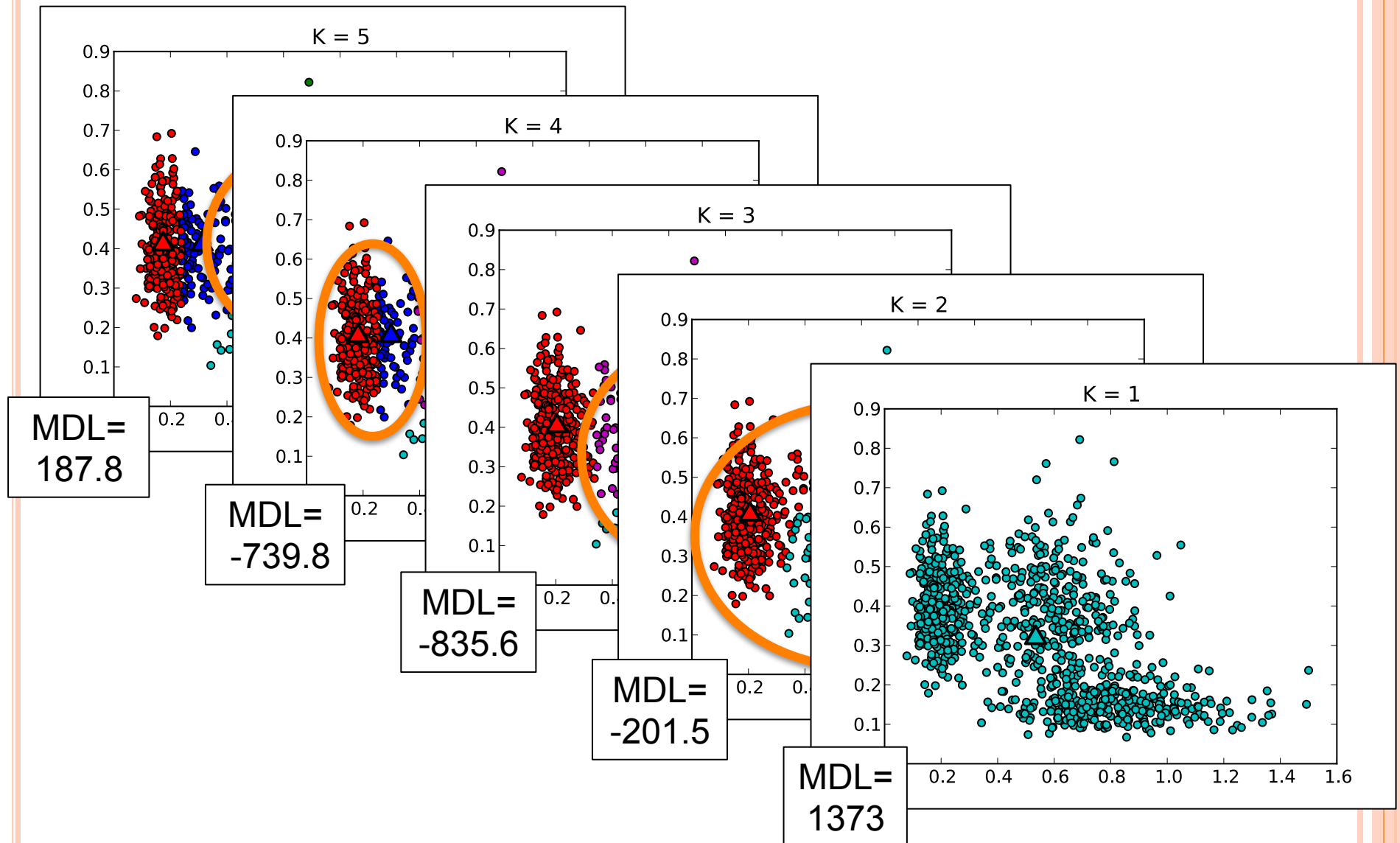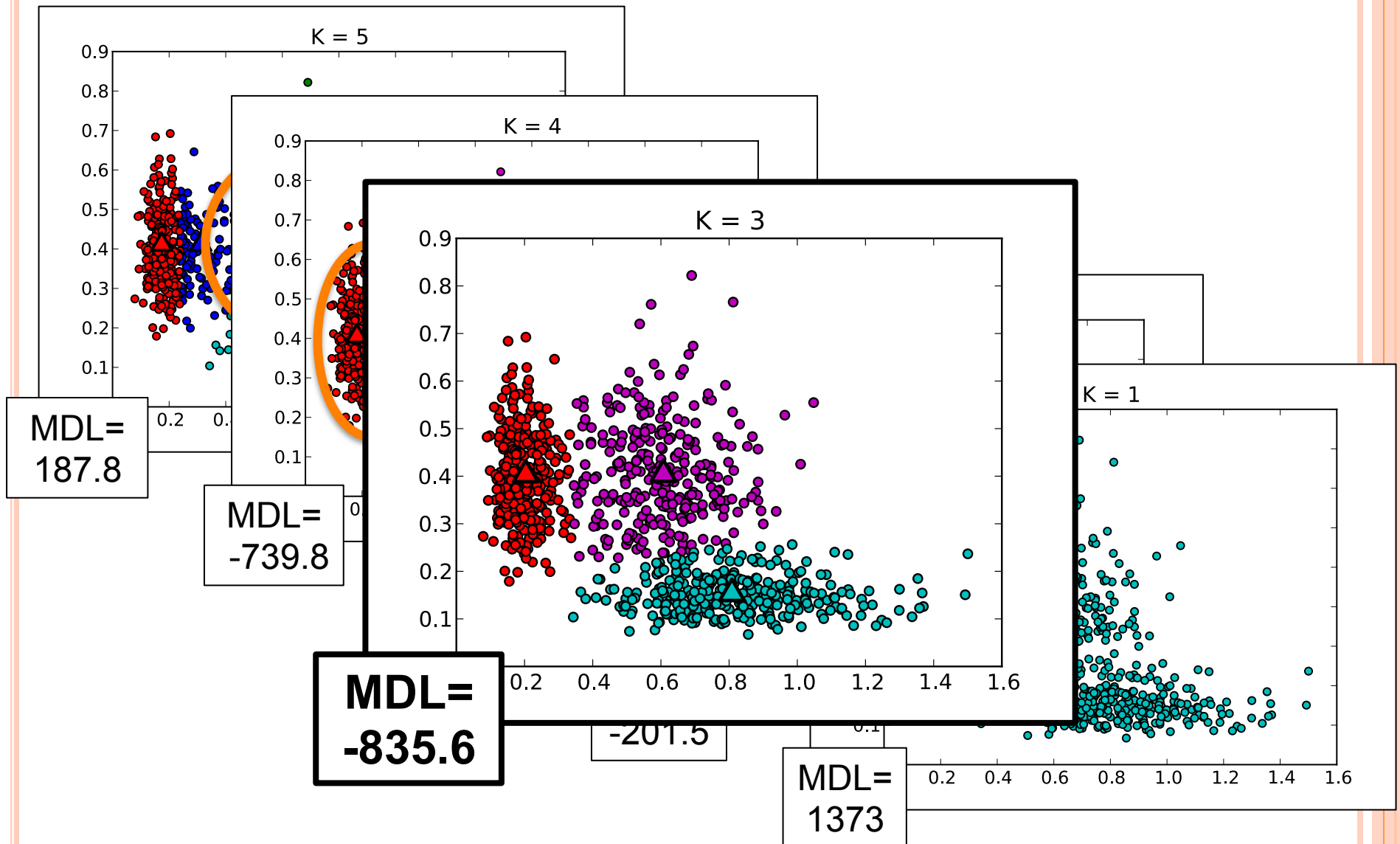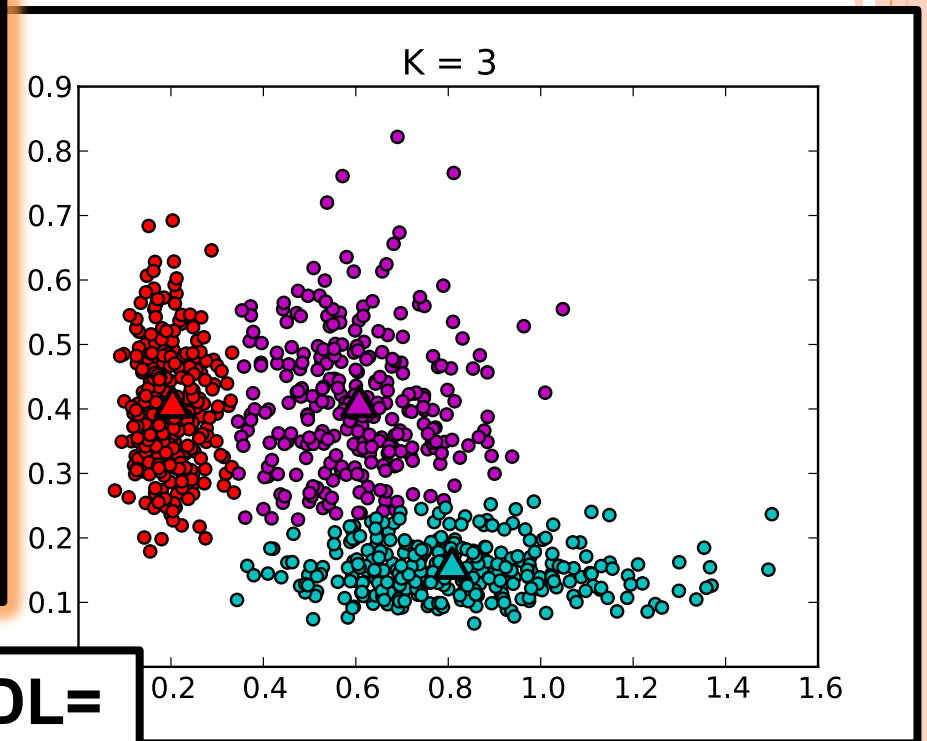# AGGLOMERATIVE CLUSTERING WITH MDL

# Agglomerative Clustering with MDL

# Agglomerative Clustering with MDL



original clusters K = 3

K = 3

MDL= -835.6

# DRUM SEPARATION SYSTEM

training data
(drum-wise audio)

performance
(raw audio)

| | |
|---|---|
| Onset Detection | Gamma Mixture Model |

cluster parameters
(drum templates)

| | |
|---|---|
| Spectrogram Slice Extraction | Non-negative Vector Decomposition |

drum activations

# DECOMPOSING ONSETS ONTO TEMPLATES

- Non-negative Vector Decomposition (NVD)
  - A simplification of Non-negative Matrix Factorization (NMF)
  - W matrix contains drum templates in its columns.

$$\min_{\vec{h}} d_{IS}(\vec{x}, W\vec{h}), \quad h_i \geq 0 \quad \forall i$$

# DECOMPOSING ONSETS ONTO TEMPLATES

o To solve this problem:

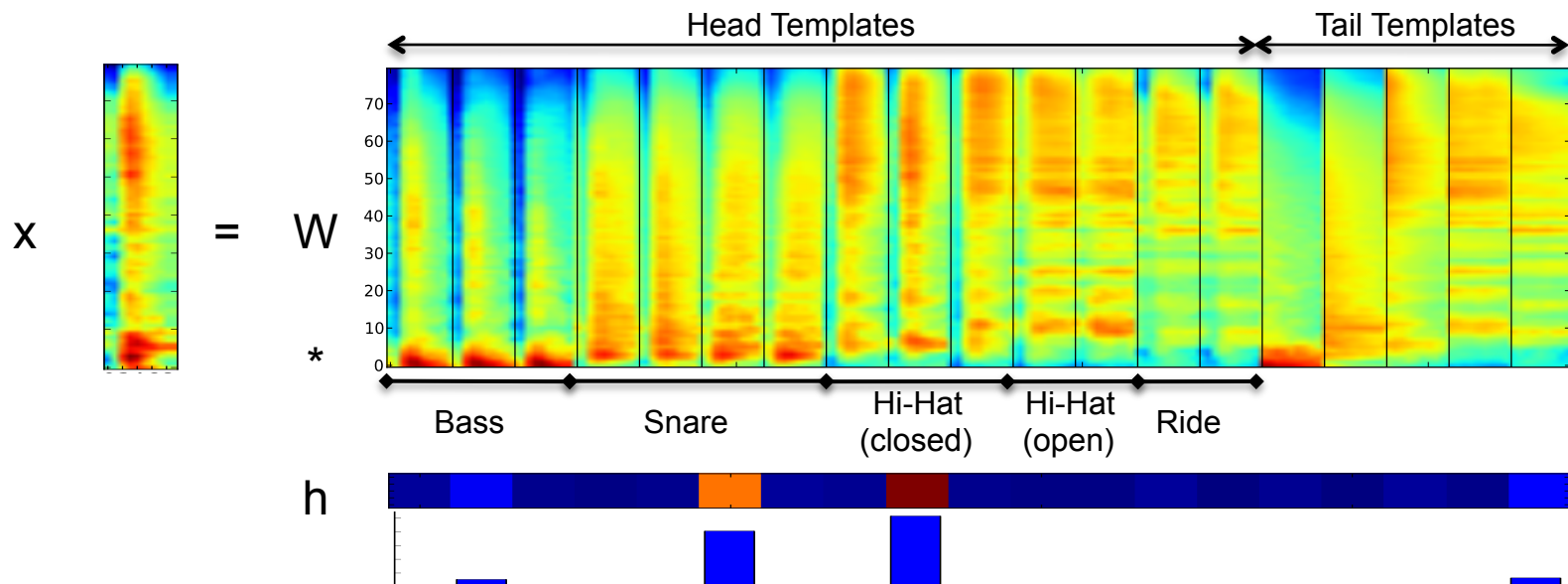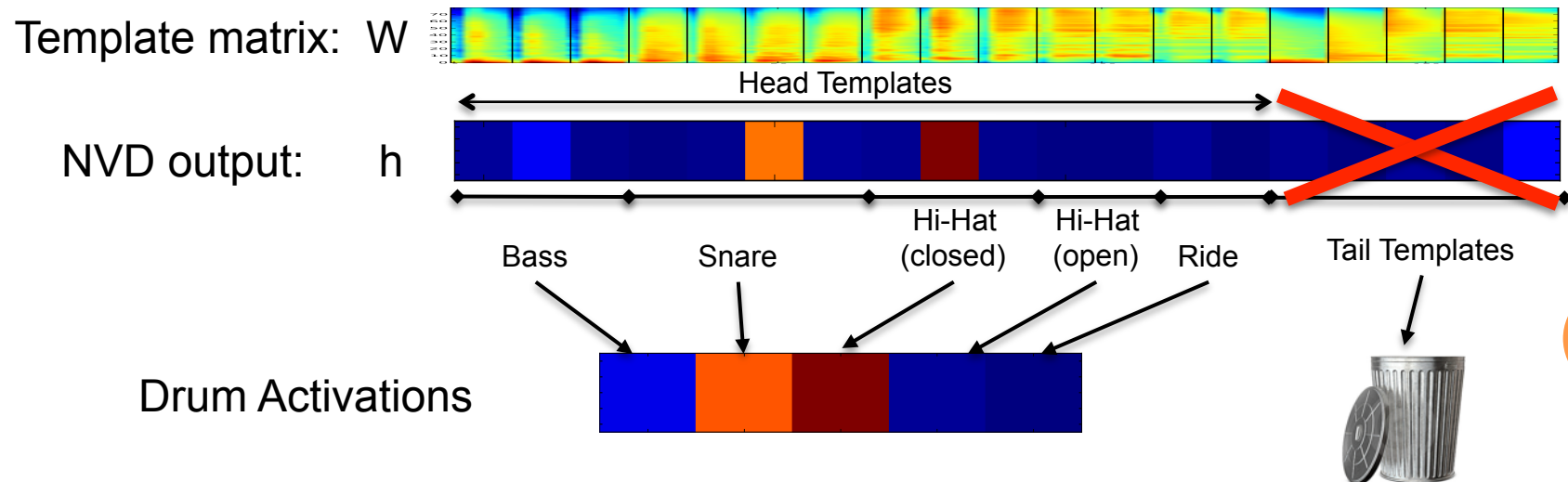$$\min_{\vec{h}} d_{IS}(\vec{x}, W\vec{h}), \quad h_i \geq 0 \quad \forall i$$

o We use the IS distance as the cost function in the above.

- While the IS distance is not strictly convex, in practice it is non-increasing under the following update rule:

$$\vec{h}_i \leftarrow \vec{h}_i \cdot \frac{W^T((W\vec{h}_i)^{\cdot -2} \cdot \vec{x}_i)}{W^T(W\vec{h}_i)^{\cdot -1}}$$

# DECOMPOSING ONSETS ONTO TEMPLATES

- What do we do with the output of NVD?
  - The **head** template activations for a single drum are **summed** to get the total activation of that drum.
  - The **tail** template activations are **discarded**.
    - They simply serve as "decoys" so that the long decay of a previous onset does not affect the current decomposition as drastically.

Template matrix: W

Head Templates

NVD output: h

Hi-Hat (closed)    Hi-Hat (open)    Ride

Bass          Snare                              Tail Templates

Drum Activations

# DRUM SEPARATION SYSTEM

training data
(drum-wise audio)

performance
(raw audio)

| Onset Detection | | Gamma Mixture Model |
|---|---|---|

cluster parameters
(drum templates)

| Spectrogram Slice Extraction | | Non-negative Vector Decomposition |
|---|---|---|

drum
activations

# BUILDING/TESTING THE SYSTEM



- Implemented in Python with Scipy
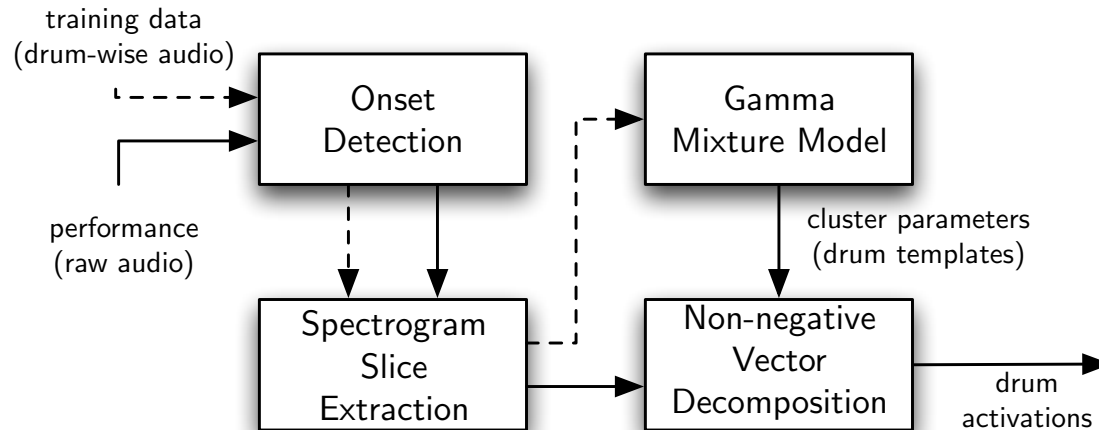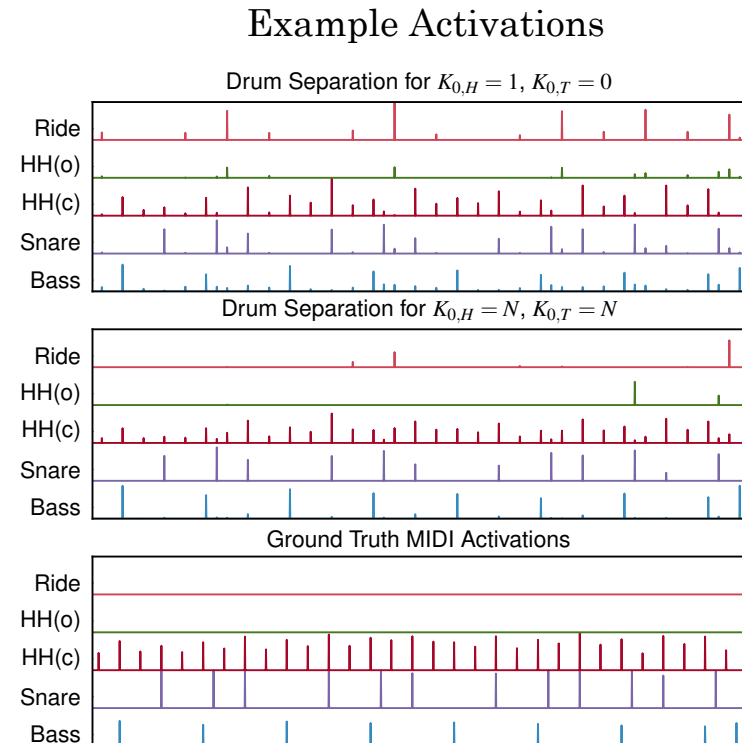  - NVD can easily be done in real-time (100ms latency)
  - Agglomerative Gamma Mixture Model training takes ~20 seconds for 5 drums.
    - *Could be reduced to < 1 sec using a GPU implementation.*
- Parameters to vary for testing:
  - Number head/tail templates per drum
    - {0, 1, MDL-optimal}
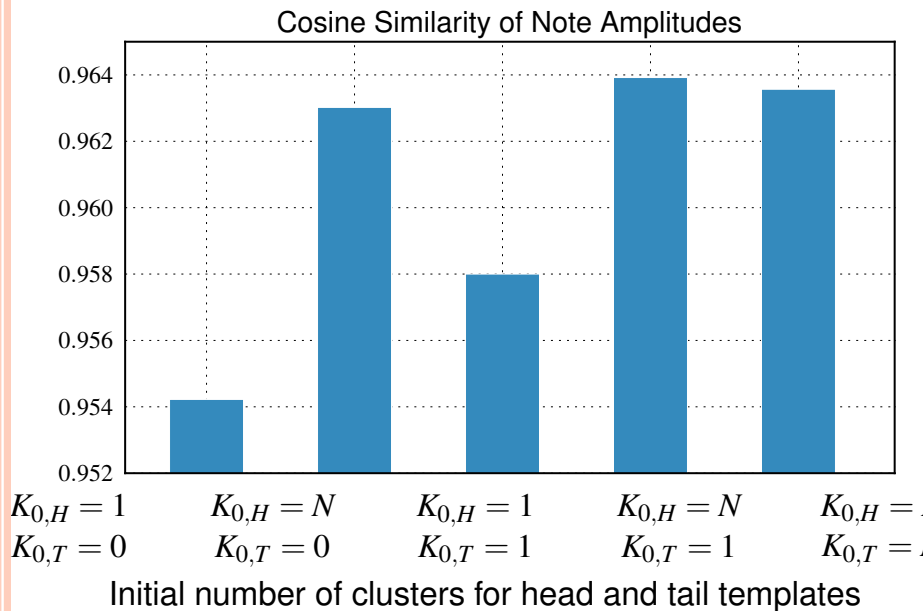
# QUANTITATIVE RESULTS

- We test using a total of 10 drum performances:
  - 10 minutes total, 2922 drum onsets
  - Recorded as midi data
    - Roland V-Drums
  - Audio created using multi-sampled drum kit
    - Superior Drummer 2.0
- Onset detection results
  - **85% recall, 99.9% precision**
- Decomposition results
  - **Cosine similarity for true activations**
  - **Amplitude sum for false activations**

Example Activations

Drum Separation for $K_{0,H} = 1, K_{0,T} = 0$

Ride | HH(o) | HH(c) | Snare | Bass

Drum Separation for $K_{0,H} = N, K_{0,T} = N$

Ride | HH(o) | HH(c) | Snare | Bass

Ground Truth MIDI Activations

Ride | HH(o) | HH(c) | Snare | Bass

$$S_{\cos}(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\| \|\vec{y}\|}$$

33

# QUANTITATIVE RESULTS

## Cosine similarity of True Activations



## Amplitude sum of False Activations



Significant improvements seen
with
- > head template
- il templates

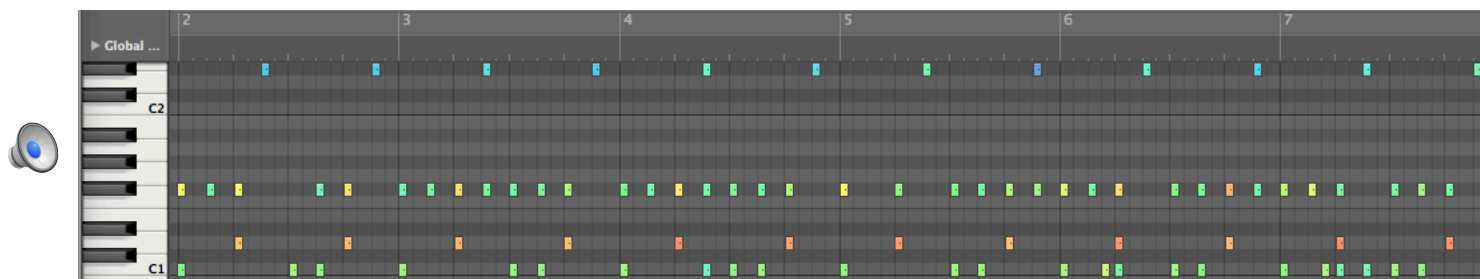# AUDIO EXAMPLES

○ Track 1 - Basic 4/4 rock beat (quantized)

Original Performance

KH=MDL-Optimal, KT=1

KH=1, KT=0

35

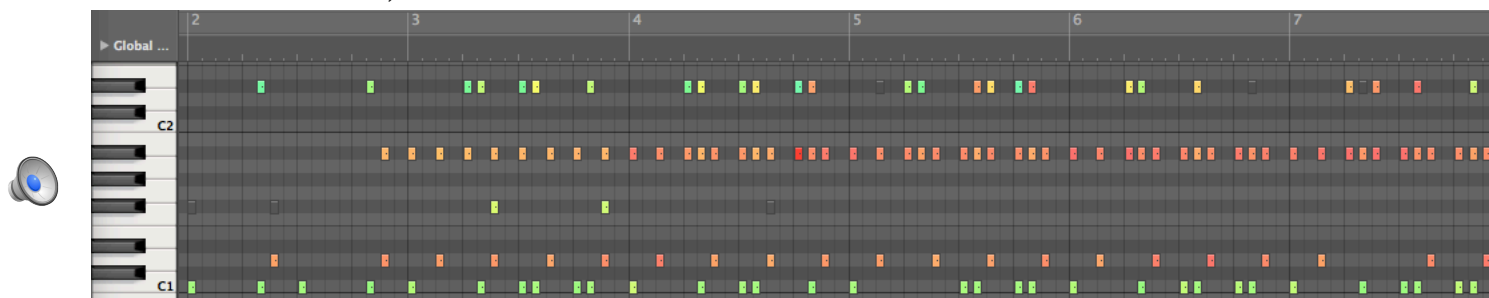# AUDIO EXAMPLES

- Track 3 - Cut time rock with open hi-hat

Original Performance
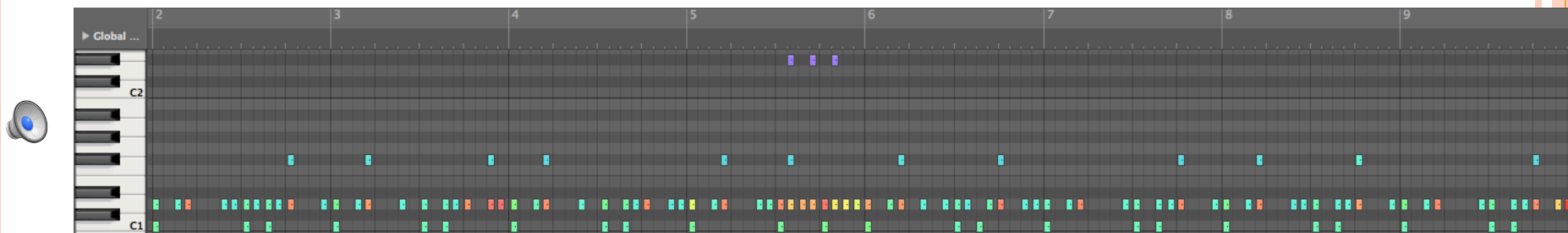


KH=MDL-Optimal, KT=1



KH=1, KT=0



36

# AUDIO EXAMPLES

- Track 7 - Accented snare drum roll.

Original Performance

KH=MDL-Optimal, KT=1

KH=1, KT=0

# SUMMARY

- Drum separation front end for a complete drum understanding system.
- Gamma Mixture Model
  - <u>Cheaper to train</u> than GMM (no covariance matrix)
  - <u>More stable</u> than GMM (no covariance matrix)
  - Allows clustering with perceptual <u>Itakura-Saito distance</u>
- Non-negative Vector Decomposition
  - Greatly improved with <u>tail templates</u> and <u>multiple head templates</u> per drum.
- Next steps
  - Explore online training of templates.
  - Integration with complete drum understanding system.

38

# Kiitos

# EXTRA SLIDES

# GAMMA MIXTURE MODEL

- Multivariate Gamma (independent components):

$$p(\vec{y}|\vec{\lambda}, k) \;=\; \prod_{i=1}^{M} \frac{\lambda_i^k y_i^{k-1} e^{-\lambda_i y_i}}{\Gamma(k)}$$

- Mixture density:

$$p(\vec{y_n}|\theta) \;=\; \sum_{l=1}^{K} \pi_l p(\vec{y_n}|\vec{\lambda}_l, k) \qquad \theta = \{\vec{\lambda}_l, \pi_l\}_{l=1}^{K}$$

$$\pi_l \;=\; p(x_n = l)$$

# THE EM ALGORITHM: GAMMA EDITION

- E-step: (compute posteriors)

$$p(x_n = l|\vec{y}_n, \theta^{(t)}) \quad = \quad \frac{\pi_l \exp\left(-k\, d_{\text{IS}}(\vec{y}_n, \vec{\mu}_l)\right)}{\sum_{j=1}^{K} \pi_j \exp\left(-k\, d_{\text{IS}}(\vec{y}_n, \vec{\mu}_j)\right)}$$

- M-step: (update parameters)

$$N_l^* \quad = \quad \sum_{n=1}^{N} p(x_n = l|\vec{y}_n, \theta^{(t)})$$

$$\vec{\lambda}_l \quad \leftarrow \quad \frac{k N_l^*}{\sum_{n=1}^{N} \vec{y}_n p(x_n = l|\vec{y}_n, \theta^{(t)})}$$

$$\pi_l \quad \leftarrow \quad \frac{N_l^*}{N}$$

# AGGLOMERATIVE CLUSTERING

- *How many clusters to train?*
- We use Minimum Description Length (MDL) to choose the number of clusters.
  - Negative log-likelihood
  - + penalty term for number of clusters.

$$\text{MDL}(K, \theta) = -\sum_{n=1}^{N} \log \left( \sum_{l=1}^{K} p(\vec{y}_n | \vec{\lambda}_l) \pi_l \right) + \tfrac{1}{2} L \log(NM)$$

$$L = KM + (K-1)$$

- 1. <u>Run EM</u> to convergence.
- 2. <u>Merge</u> the two most similar clusters.
- 3. <u>Repeat</u> 1,2 until we have a single cluster.
- 4. *Choose parameter set with smallest MDL.*