

TEACHING COMPUTERS TO LISTEN TO MUSIC

Eric Battenberg

<http://ericbattenberg.com>

Gracenote, Inc.

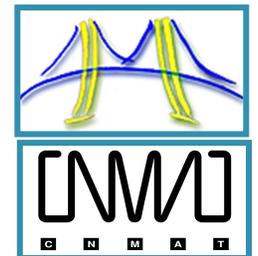


Previously:

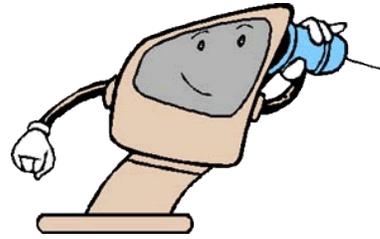
UC Berkeley, Dept. of EECS

Parallel Computing Laboratory

CNMAT (Center for New Music and Audio Technologies)



MACHINE LISTENING



- Speech processing

- Speech processing makes up the vast majority of funded machine listening research.
- Just as there's more to Computer Vision than OCR, there's more to machine listening than speech recognition!



- Audio content analysis

- Audio fingerprinting (e.g. Shazam, Gracenote)
- Audio classification (music, speech, noise, laughter, cheering)
- Audio event detection (new song, channel change, hotword)



- Content-based Music Information Retrieval (MIR)

- Today's topic



GETTING COMPUTERS TO “LISTEN” TO MUSIC

- Not trying to get computers to “listen” for enjoyment.
- More accurate: Analyzing music with computers.
- What kind of information to get out of the analysis?
 - What instruments are playing?
 - What is the mood?
 - How fast/slow is it (tempo)?
 - What does the singer sound like?
 - How can I play this song on the guitar/drums?

Ben Harper

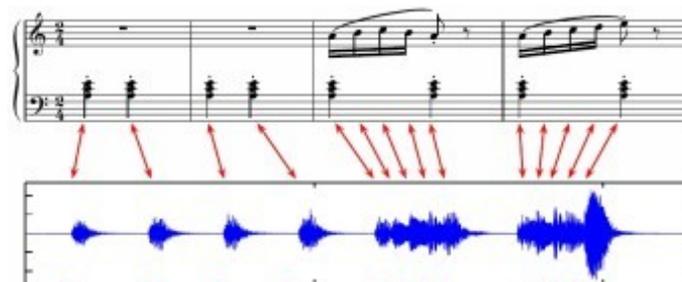
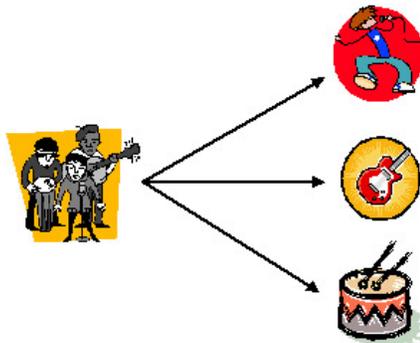


James Brown



CONTENT-BASED MUSIC INFORMATION RETRIEVAL

- Many tasks:
 - Genre, mood classification, auto-tagging
 - Beat tracking, tempo detection
 - Music similarity, playlisting, recommendation
 - Heavily aided by collaborative filtering
 - Automatic music transcription
 - Source separation (voice, drums, bass...)
 - Music segmentation (verse, chorus)



listen to the music that best suits my mood or my activity

RELAX happy calm CHILLOUT dreamy
JUST WOKE UP optimistic SUNNY DAY melancholy
SUMMER GOOD KARMA STUDYING sleepy
SPRING CLEANING angry cool LOST IN THOUGHT
energetic sad beautiful romantic digital need of love
SUNDAY MORNING crazy IT'S RAINING sweet READING sexy

Sun 15 Jan | 1:20 PM
17/42 Music/playlist
Duration 50:53

Jacksoul - Sleepless	4:32	▲	▼	×
Blaq Lily - Dreams	3:22	▲	▼	×
Frou Frou - Breathe In	4:28	▲	▼	×
dZihan & Kamien - Just You and I		▲	▼	×
Christina Smith & Jean Hewson - Liar's	3:37	▲	▼	×
Poe - Control	6:03	▲	▼	×
Kinnie Starr - Soar	2:49	▲	▼	×
Madonna - Skin	6:22	▲	▼	×
Macy Gray - Caligula	4:40	▲	▼	×
Lifehouse - Hanging By A Moment	3:35	▲	▼	×
Susheela Raman - Mamavatu	3:56	▲	▼	×
Tara MacLean - Poor Boy	3:55	▲	▼	×

TALK SUMMARY

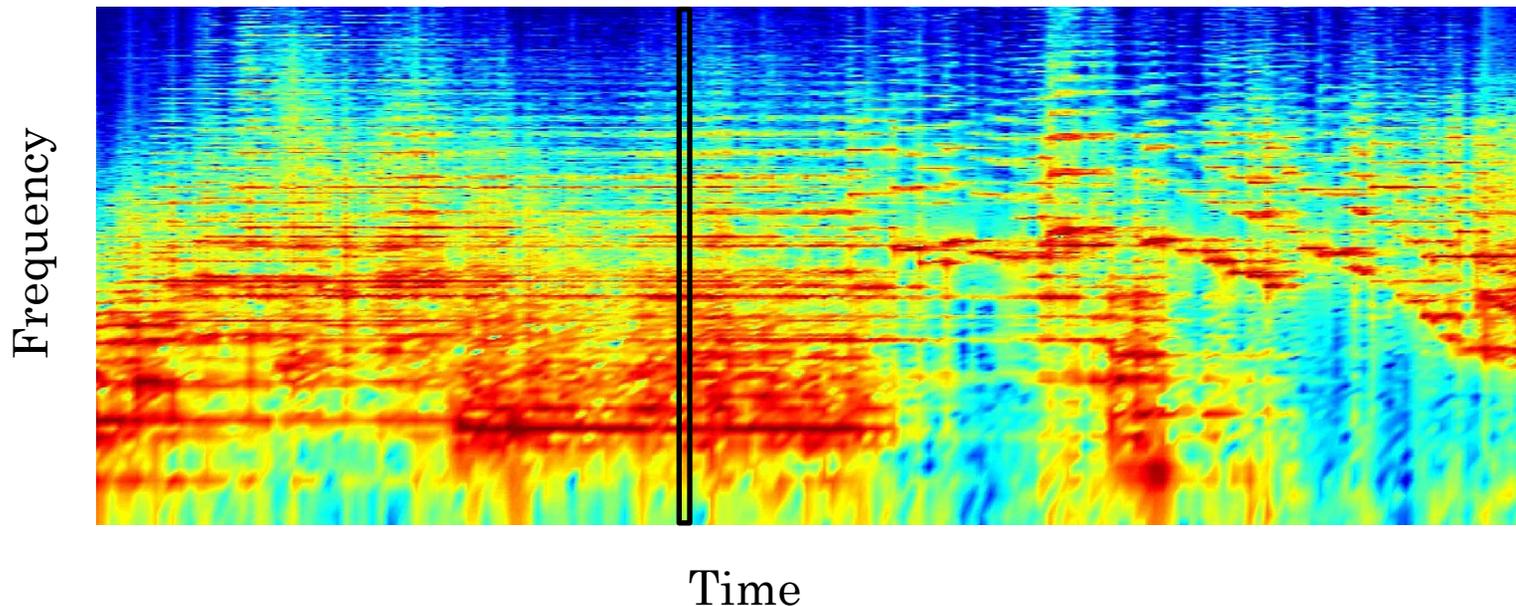
- Introduction to Music Information Retrieval
 - Some common techniques
 - Exciting new research directions
- Live Drum Understanding
 - Drum detection/transcription
 - Drum pattern analysis

TALK SUMMARY

- **Introduction to Music Information Retrieval**
 - Some common techniques
 - Exciting new research directions
- **Live Drum Understanding**
 - Drum detection/transcription
 - Drum pattern analysis

QUICK LESSON: THE SPECTROGRAM

- The spectrogram: Very common feature used in audio analysis.
- Time-frequency representation of audio.
- Take FFT of adjacent frames of audio samples, put them in a matrix.
- Each column shows frequency content at a particular time.

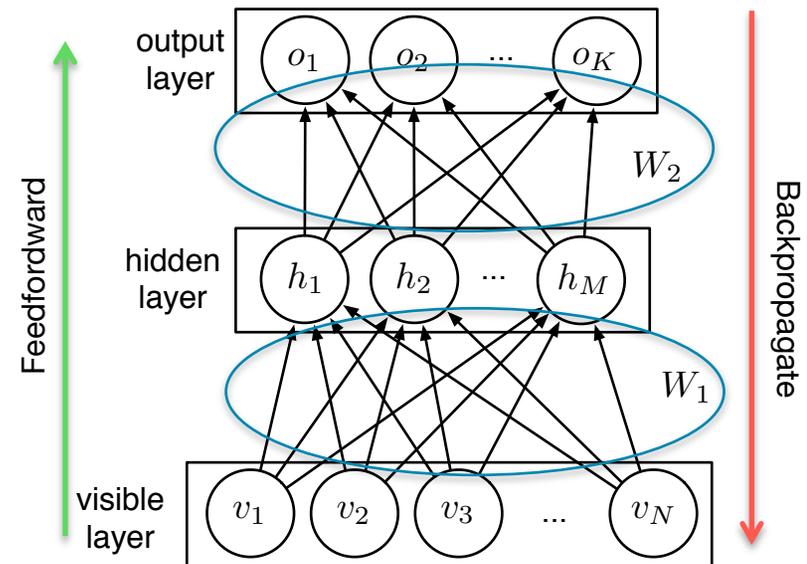
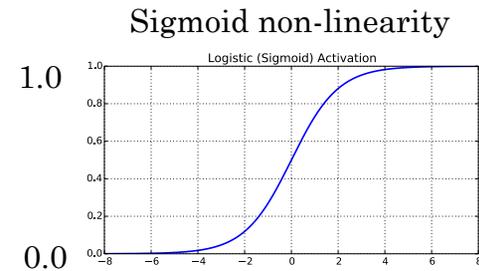


GENRE/MOOD CLASSIFICATION: TYPICAL APPROACHES

- Typical approach:
 - Extract a **bunch of hand-designed features** describing small windows of the signal (e.g., spectral centroid, kurtosis, harmonicity, percussiveness, MFCCs, 100's more...).
 - **Train a GMM or SVM** to predict genre/mood/tags by either:
 - Summarizing a song using mean/variance of each feature
 - Log-likelihood sum across frames (GMM) or frame-wise voting (SVM)
- Pros:
 - Works fairly well, was state of the art for a while.
 - Well understood models, implementations widely available.
- Cons:
 - Bag-of-frames style approach lacks ability to describe rhythm and temporal dynamics.
 - Getting further improvements requires hand designing more features.

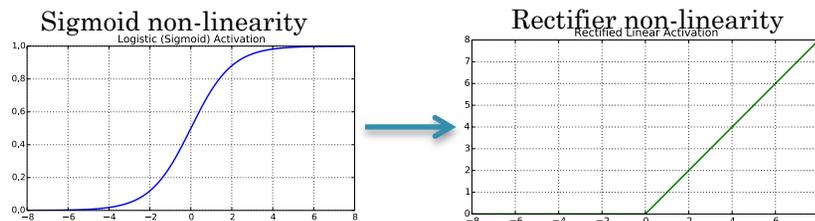
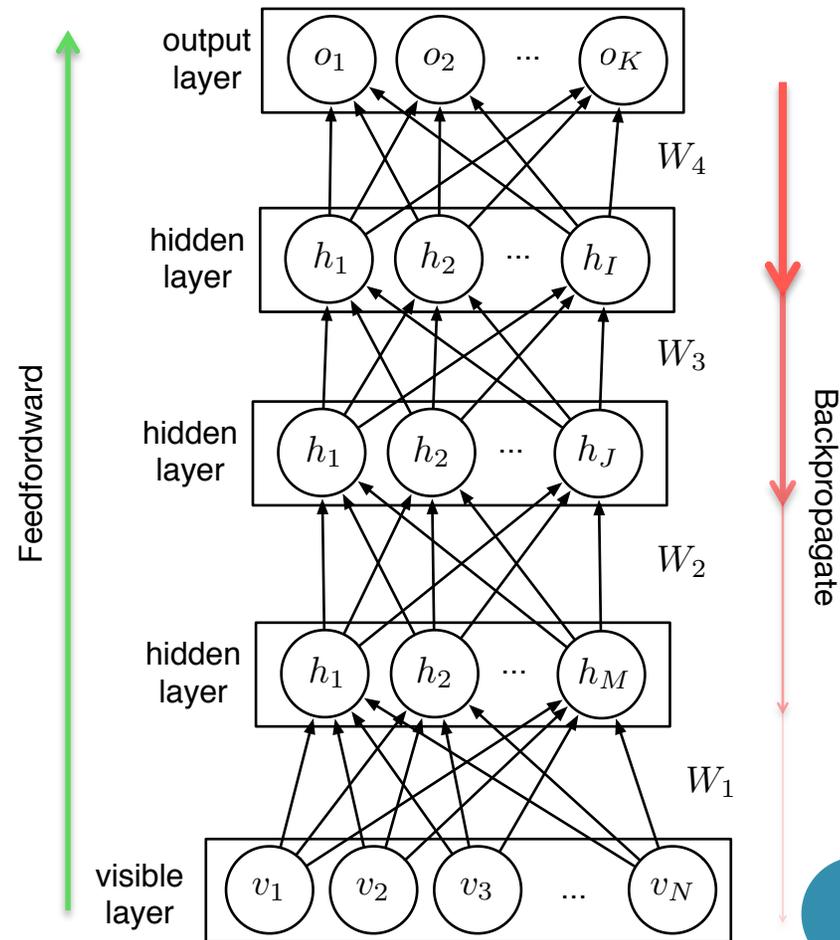
LEARNING FEATURES: NEURAL NETWORKS

- Each layer computes a non-linear transformation of the previous layer.
 - Linear transformation (weight matrices)
 - + non-linearity (e.g. sigmoid $[\sigma]$)
 - $h = \sigma(W_1v)$, $o = \sigma(W_2h)$
- Train to minimize output error.
- Each hidden layer can be thought of as a set of features.
- Train using backpropagation.
- Iterative steps:
 - Compute activations 
 - Compute output error
 - Backprop. error signal 
 - Compute gradients
 - Update all weights. 
- Resurgence of neural networks:
 - More compute
 - More data
 - A few new tricks...



DEEP NEURAL NETWORKS

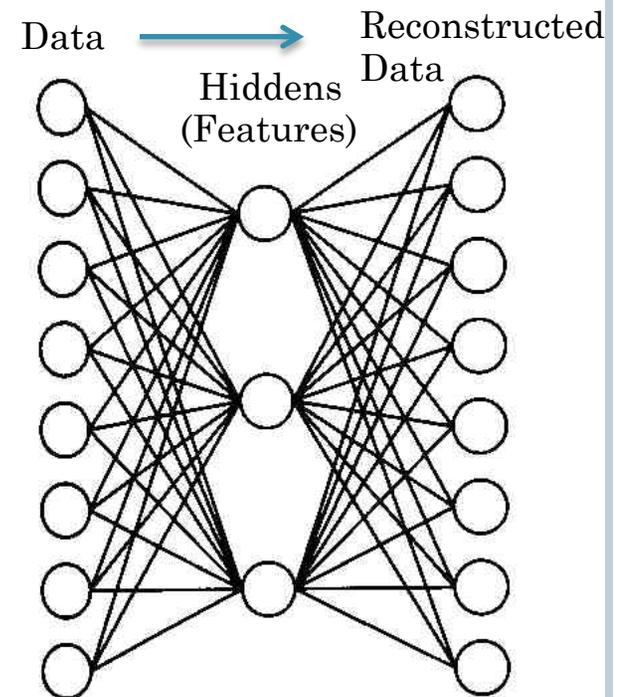
- Deep Neural Networks
 - Millions to billions of parameters
 - Many layers of “features”
 - Achieving state of the art performance in vision and speech tasks.
- Problem: Vanishing error signal.
 - Weights of lower layers do not change much.
- Solutions:
 - Train for a really long time. ☹
 - Pre-train each hidden layer as an autoencoder. [Hinton, 2006]
 - Rectified Linear Units [Krizhevsky, 2012]



AUTOENCODERS AND UNSUPERVISED FEATURE LEARNING

- Many ways to learn features in an unsupervised way:
- Autoencoders – train a network to reconstruct the input
 - Restricted Boltzmann Machine (RBM) [Hinton, 2006]
 - Denoising Autoencoders [Vincent, 2008]
 - Sparse Autoencoders
- Clustering – K-Means, mixture models, etc.
- Sparse Coding – learn overcomplete dictionary of features with sparsity constraint

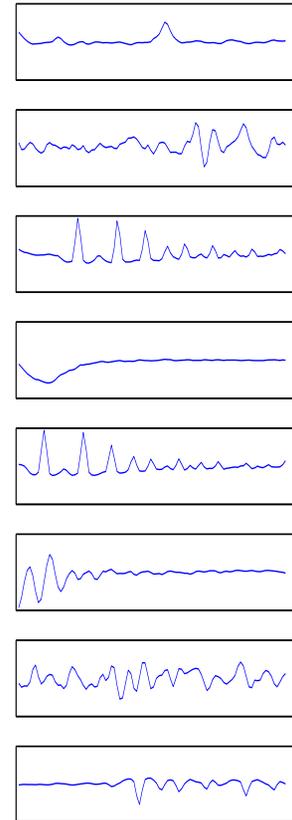
Autoencoder:
Train to reconstruct input



GENRE/MOOD CLASSIFICATION: NEWER APPROACHES

- Newer approaches to feature extraction:
 - Learn spectral features using Restricted Boltzmann Machines (RBMs) and Deep Neural Networks (DNN) [Hamel, 2010] – good genre performance.
 - Learn sparse features using Predictive Sparse Decomposition (PSD) [Henaff, 2011] – good genre performance
 - Learn beat-synchronous rhythm and timbre features with RBMs and DNNs [Schmidt, 2013] – improved mood performance
 - Tune multi-layer wavelet features called Deep Scattering Spectrum [Anden, 2013]* - state-of-the-art genre performance
- Pros:
 - Hand-designing individual features is not required.
 - Computers can learn complex high-order features that humans cannot hand code.
- Further work:
 - More work on incorporating context, rhythm, and temporal dynamics into feature learning

Some learned
frequency features

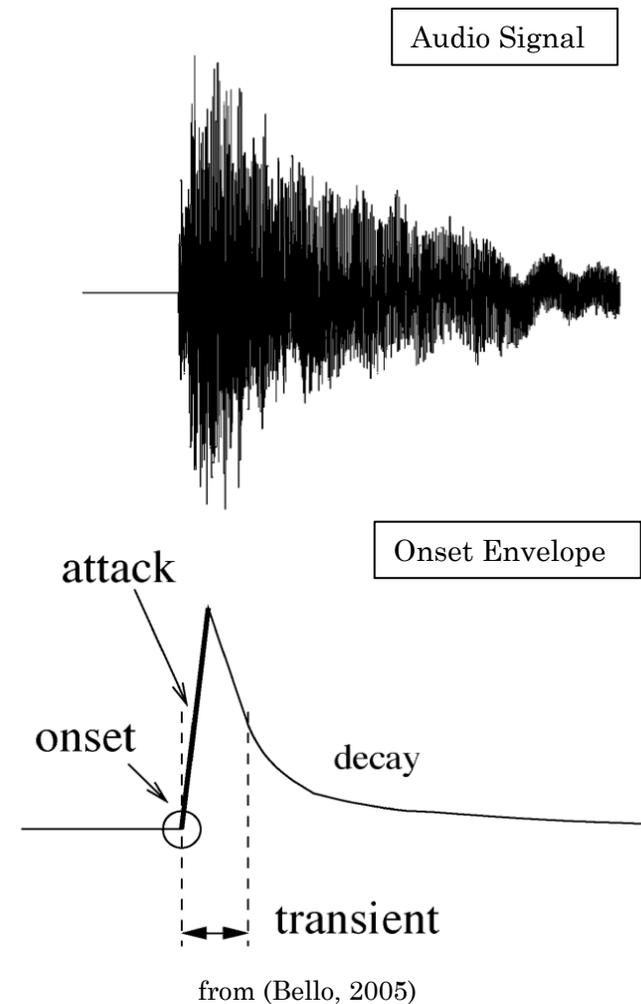


[Henaff, 2011]

* Current state-of-the art on GTZAN genre dataset

ONSET DETECTION

- Onset detection is important for music transcription, beat tracking, tempo detection, and rhythm summarization.
- Describing an onset:
 - Transient
 - Short interval when music signal evolves unpredictably.
 - Attack
 - Amplitude envelope increasing.
 - Decay
 - Amplitude envelope decreasing.
 - Onset
 - Point in time chosen to represent beginning of transient.
- Onset detection can be hard for certain instruments with ambiguous attacks or when a note changes without a new attack (legato).



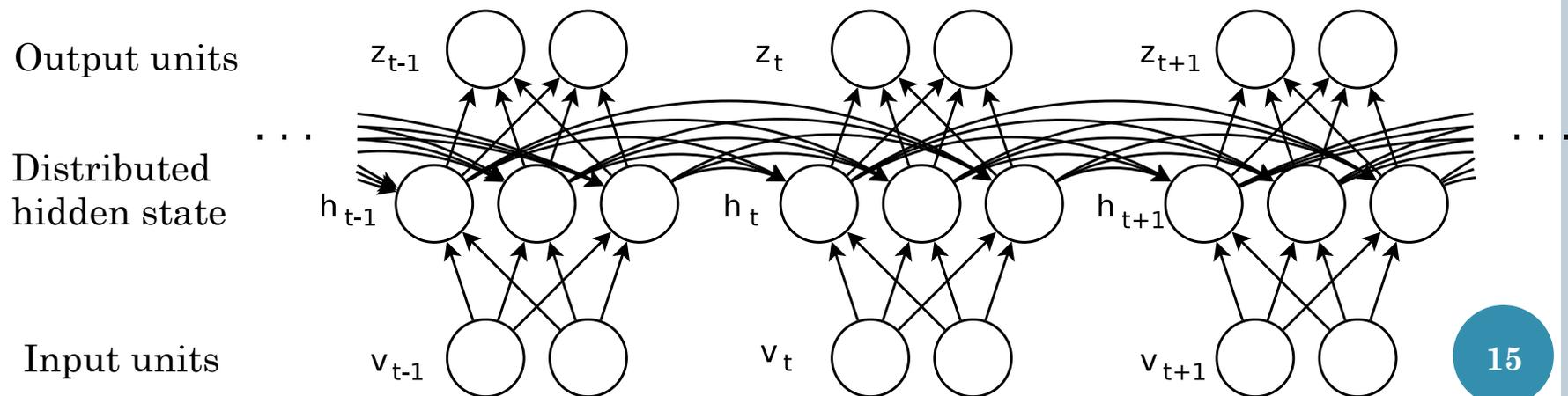
ONSET DETECTION: TYPICAL APPROACHES

- Computing an Onset Detection Function (ODF):
 - Derivative of energy envelope
 - Derivative of band-wise log-energies [Klapuri, 2006]
 - Complex spectral domain (difference with predictions of phase and magnitude) [Bello, 2004]
- Choosing onsets using the ODF:
 - Peak pick local maxima above a dynamic threshold [Bello, 2005].
- Pros:
 - Simple to implement.
 - Works fairly well (60-80% accurate)
- Cons
 - Lots of hand tuning of thresholds
 - No machine learning

RECURRENT NEURAL NETWORKS

- Non-linear sequence model
- Hidden units have connections to previous time step
- Unlike HMMs, can model **long-term dependencies** using distributed hidden state.
- Recent developments (+ more compute) have made them much more feasible to train.

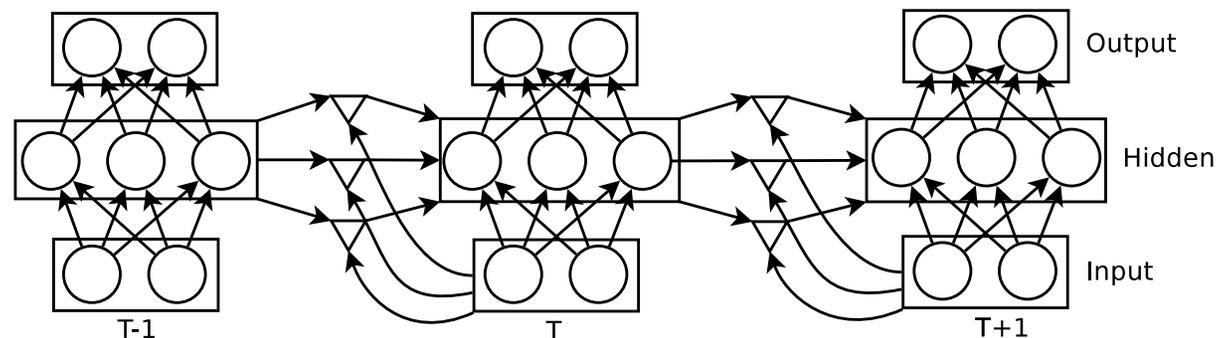
Standard Recurrent Neural Network



from [Sutskever, 2013]

TRAINING AN RNN ON WIKIPEDIA

- Train RNN to predict next character (not word)
- Multiplicative RNN [Sutskever, 2013]

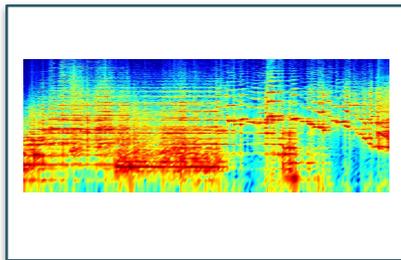


- Text generation demo: <http://www.cs.toronto.edu/~ilya/rnn.html>
- *The machine learning meetup in San Francisco is considered enormously emphasised. While as a consequence of these messages are allocated in the environment to see the ideas and pollentium changes possible with the Machinese gamma integrals increase, then the prefix is absent by a variety of fresh deeperwater or matter level on 2 and 14, yet the...*

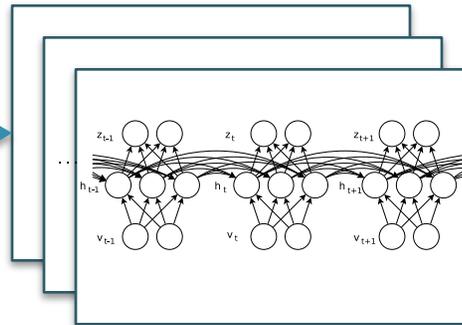
ONSET DETECTION: STATE-OF-THE-ART

- Using Recurrent Neural Networks (RNN) [Eyben, 2010], [Böck, 2012]

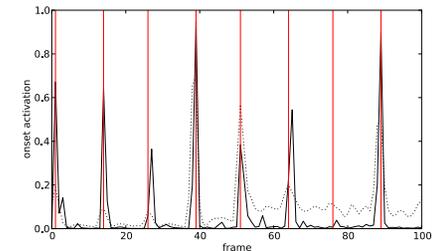
Spectrogram + Δ 's
@ 3 time-scales



3 layer RNN



Onset Detection
Function



- RNN output trained to predict onset locations.
- 80-90% accurate
- Can improve with more labeled training data, or possibly more unsupervised training.

OTHER EXAMPLES OF RNNs IN MUSIC

- Examples:
 - Blues improvisation (with LSTM RNNs) [Eck, 2002]
 - Polyphonic piano note transcription [Böck, 2012]
 - Classical music generation and transcription [Boulanger-Lewandowski, 2012]
- Other distributed-state sequence models include:
 - Recurrent Temporal Restricted Boltzmann Machine [Sutskever, 2013]
 - Conditional Restricted Boltzmann Machine [Taylor, 2011]
 - Used in drum pattern analysis later.
- RNNs are a promising way to model longer-term contextual and temporal dependencies present in music.

TALK SUMMARY

○ Introduction to MIR

- Common techniques:
 - Genre/Mood Classification, Onset Detection
- Exciting new research directions:
 - Feature learning, temporal models (RNNs)

○ Live Drum Understanding

- Drum detection/transcription
- Drum pattern analysis

TALK SUMMARY

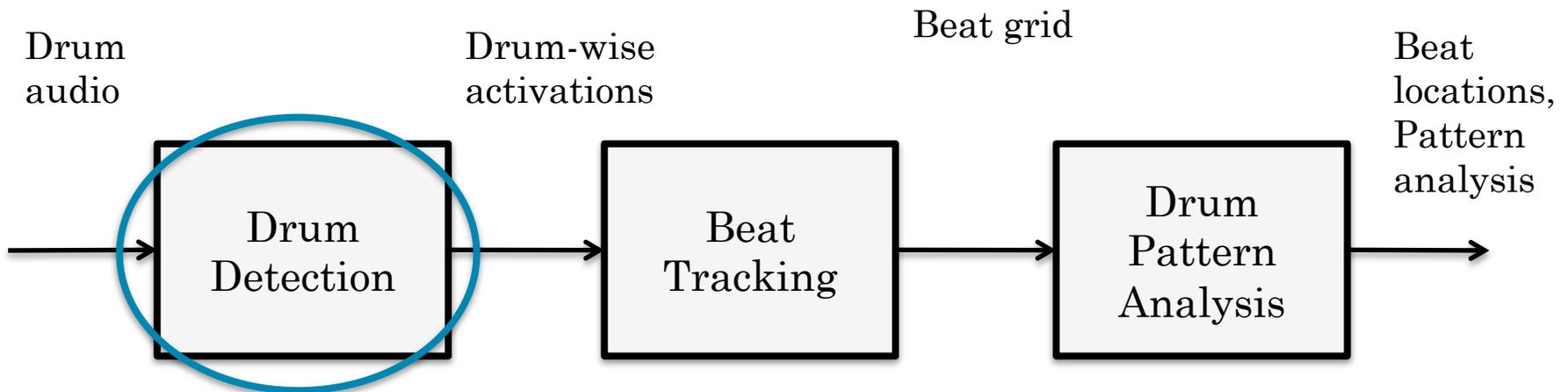
- Introduction to MIR
 - Common techniques:
 - Genre/Mood Classification, Onset Detection
 - Exciting new research directions:
 - Feature learning, temporal models (RNNs)
- **Live Drum Understanding**
 - Drum detection/transcription
 - Drum pattern analysis

TOWARD COMPREHENSIVE RHYTHMIC UNDERSTANDING

- Or “Live Drum Understanding”
- Goal: Go beyond simple beat tracking to 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



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

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

REQUIREMENTS FOR DRUM DETECTION

- Real-Time/Live operation
- Useful with any percussion setup.
 - Before a performance, we can quickly train the system for a particular percussion setup.
- Amplitude (dynamics) information.

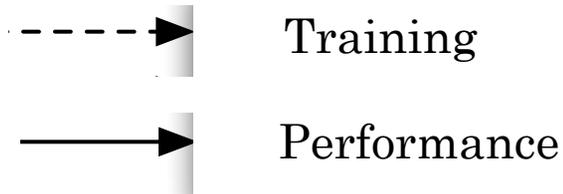
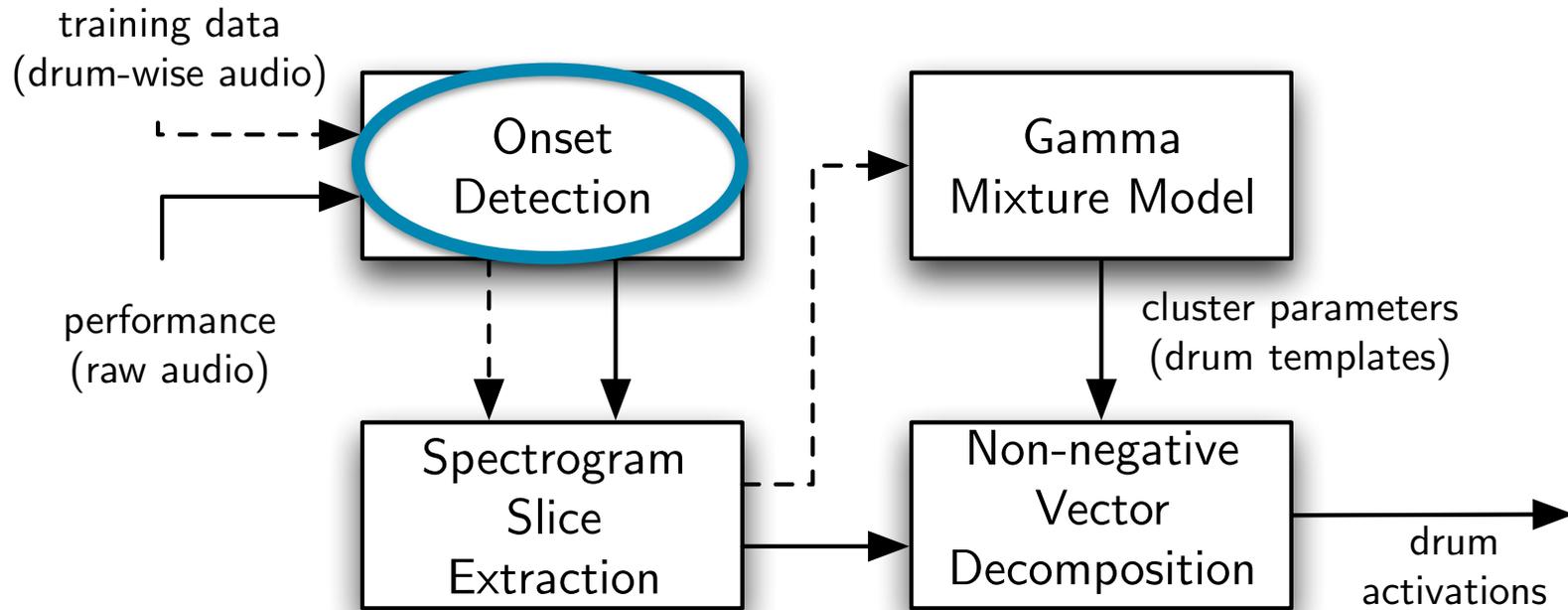


DRUM DETECTION: MAIN POINTS

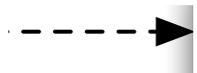
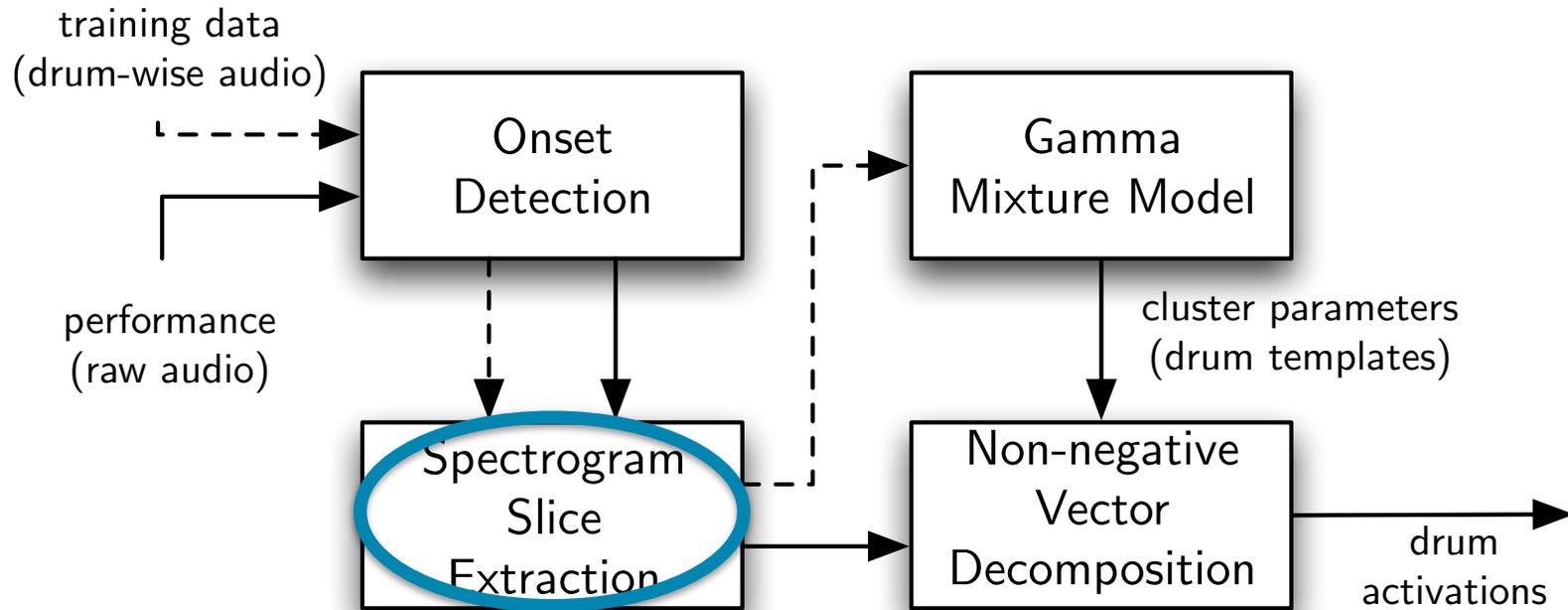


- Gamma Mixture Model
 - *For learning spectral drum templates.*
 - **Cheaper to train** than GMM
 - **More stable** than GMM
- Non-negative Vector Decomposition (NVD)
 - *For computing template activations from drum onsets.*
 - Learning **multiple templates per drum** improves separation.
 - The use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

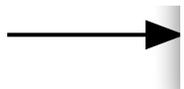
DRUM DETECTION SYSTEM



DRUM DETECTION SYSTEM



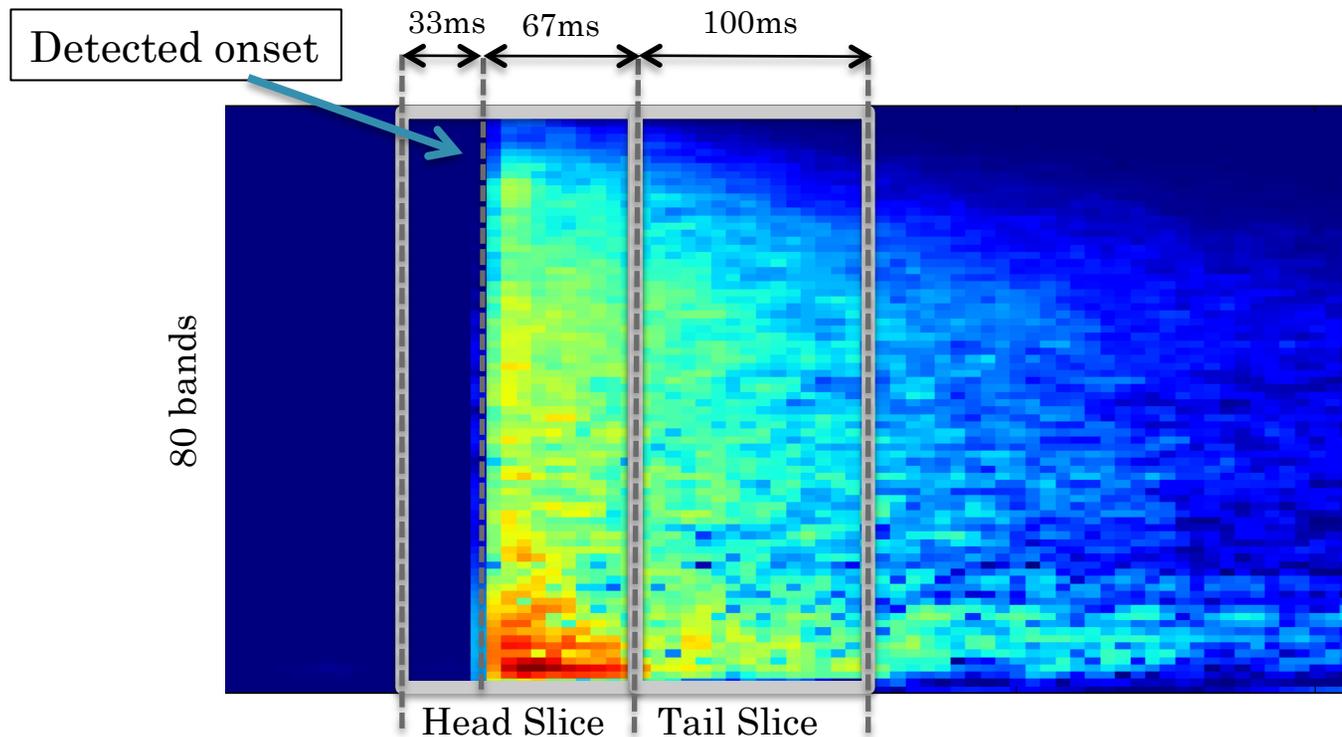
Training



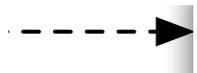
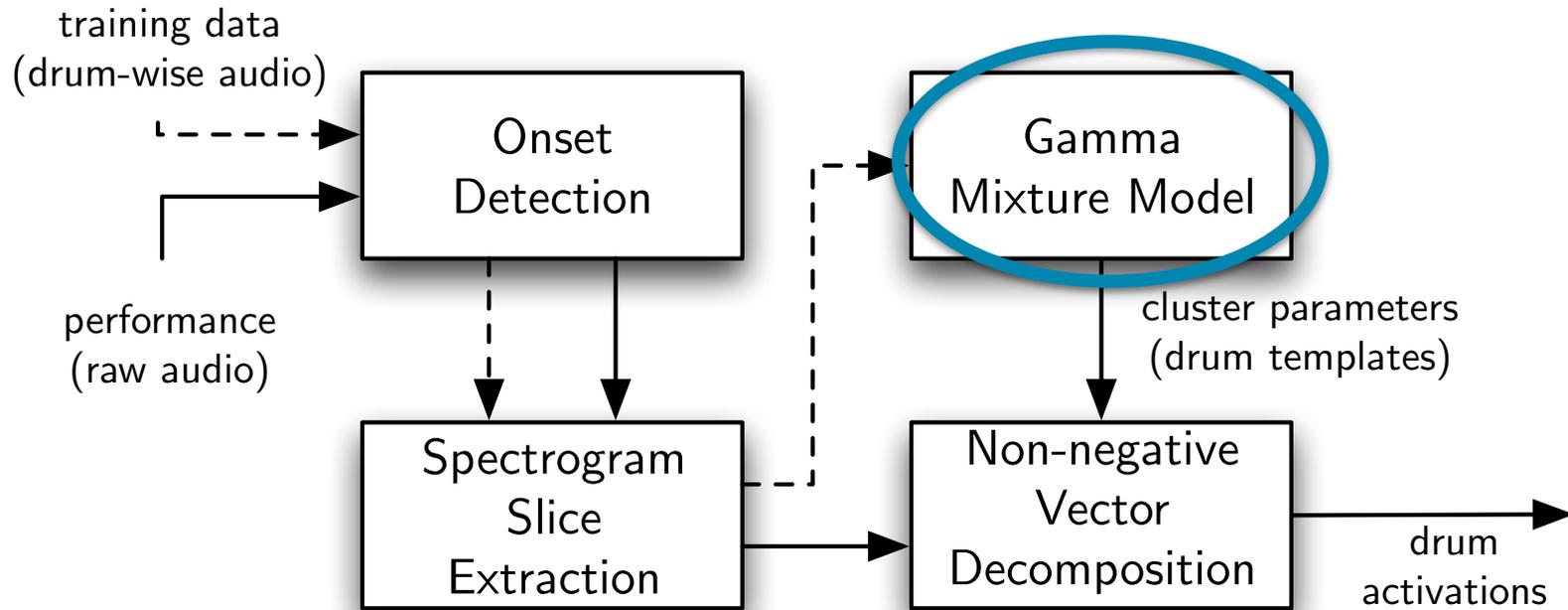
Performance

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.



DRUM DETECTION SYSTEM



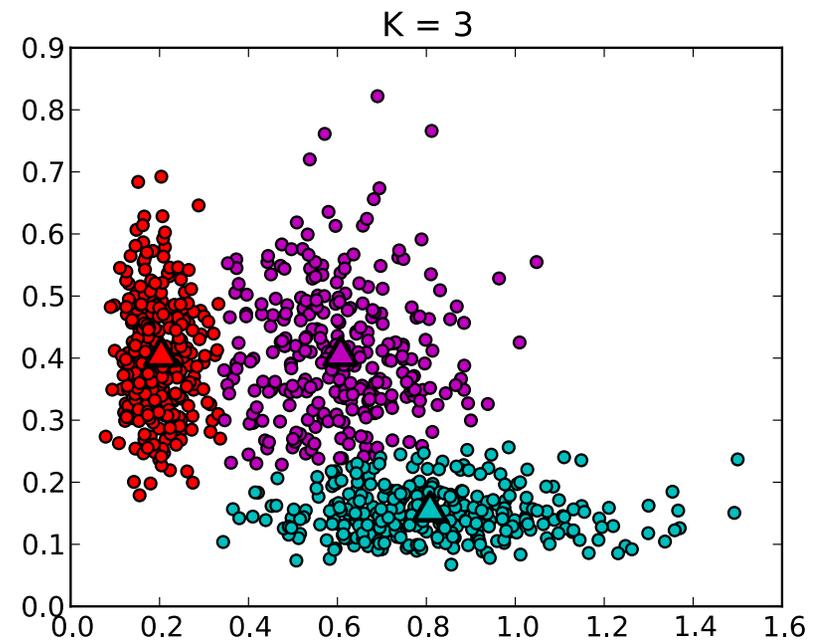
Training



Performance

TRAINING DRUM TEMPLATES

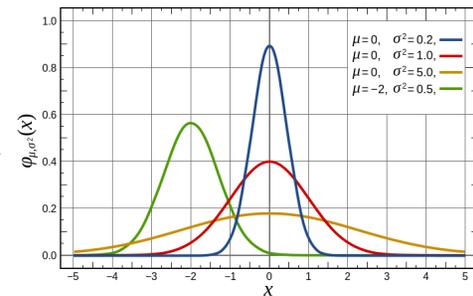
- Instead of taking an “average” of all training slices for a single drum...
- Cluster them and use the cluster centers as the drum templates.
 - This gives us multiple templates per drum...
 - Which helps represent the variety of sounds that can be made by a single drum.



CLUSTERING USING MIXTURE MODELS

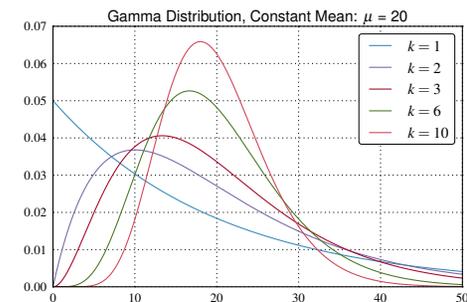
- Train using the Expectation-Maximization (EM) algorithm.
- Gaussian Mixture Model (GMM)
 - Covariance matrix – expensive to compute, possibly unstable when data is lacking.
 - Enforces a (scaled,squared) Euclidean distance measure.

$$d_{Euc}(x, y) = \sum_i (x_i - y_i)^2$$



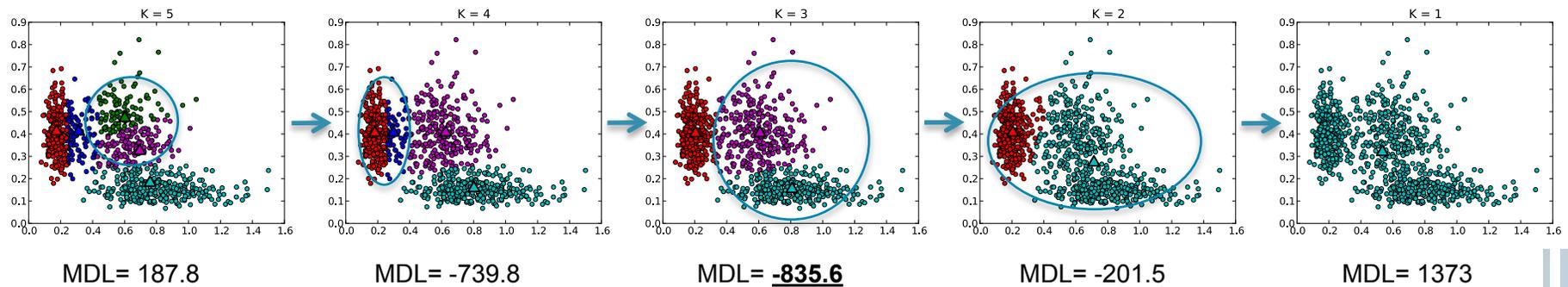
- Gamma Mixture Model
 - Single mean vector per component
 - Variance increases with mean (like human hearing).
 - Enforces an Itakura-Saito (IS) divergence measure
 - A scale-invariant perceptual distance between audio spectra.

$$d_{IS}(x, y) = \sum_i \frac{x_i}{y_i} - \log \frac{x_i}{y_i} - 1$$



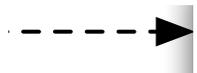
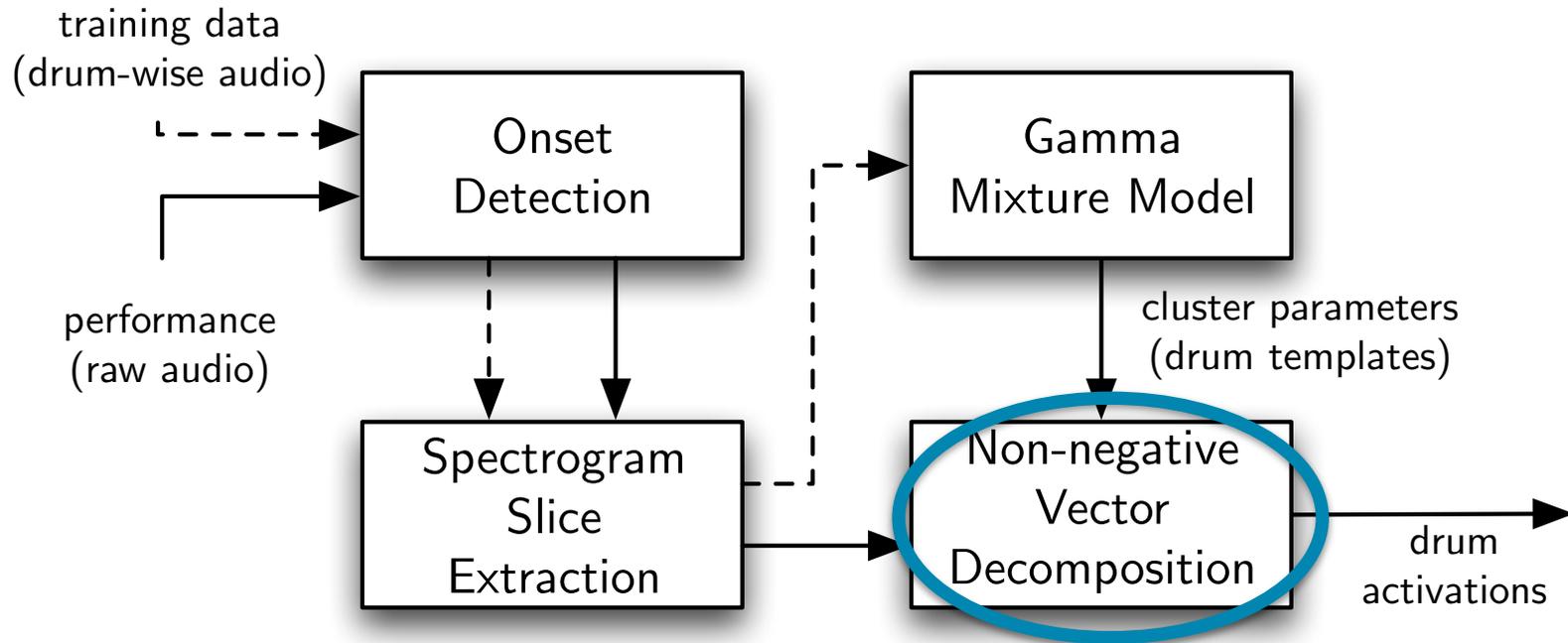
AGGLOMERATIVE CLUSTERING

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



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

DRUM DETECTION SYSTEM



Training

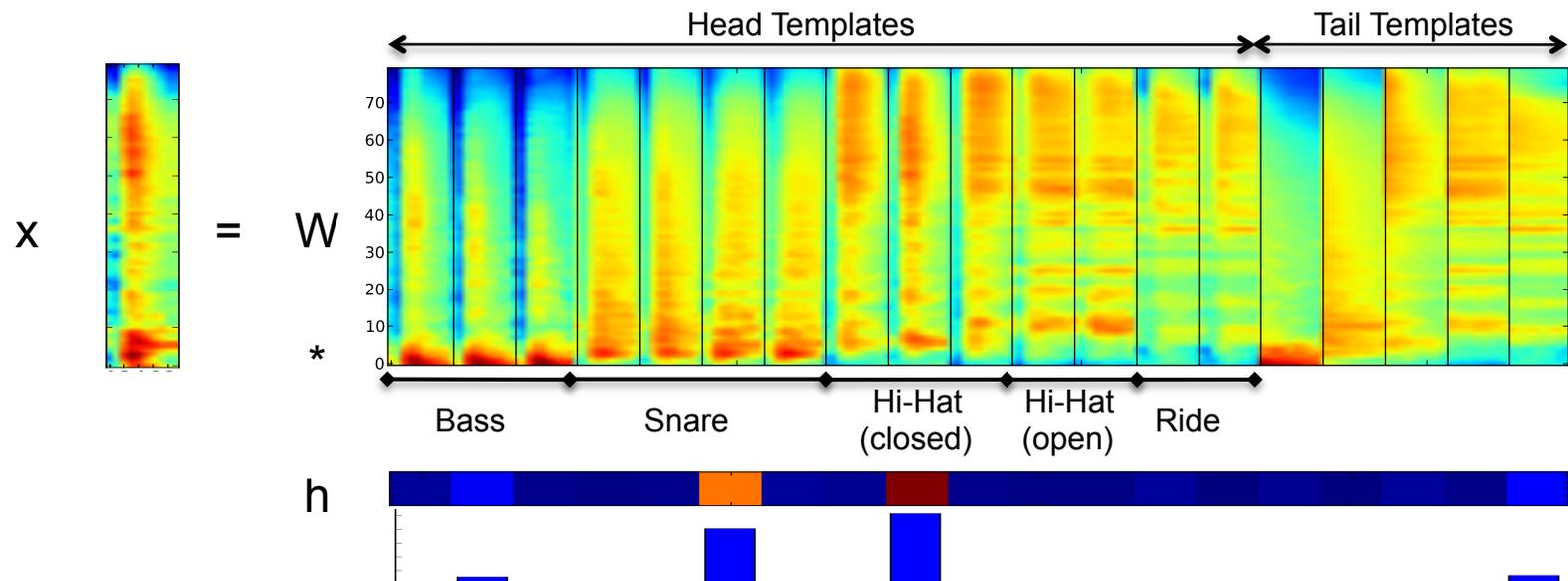


Performance

DECOMPOSING ONSETS ONTO TEMPLATES

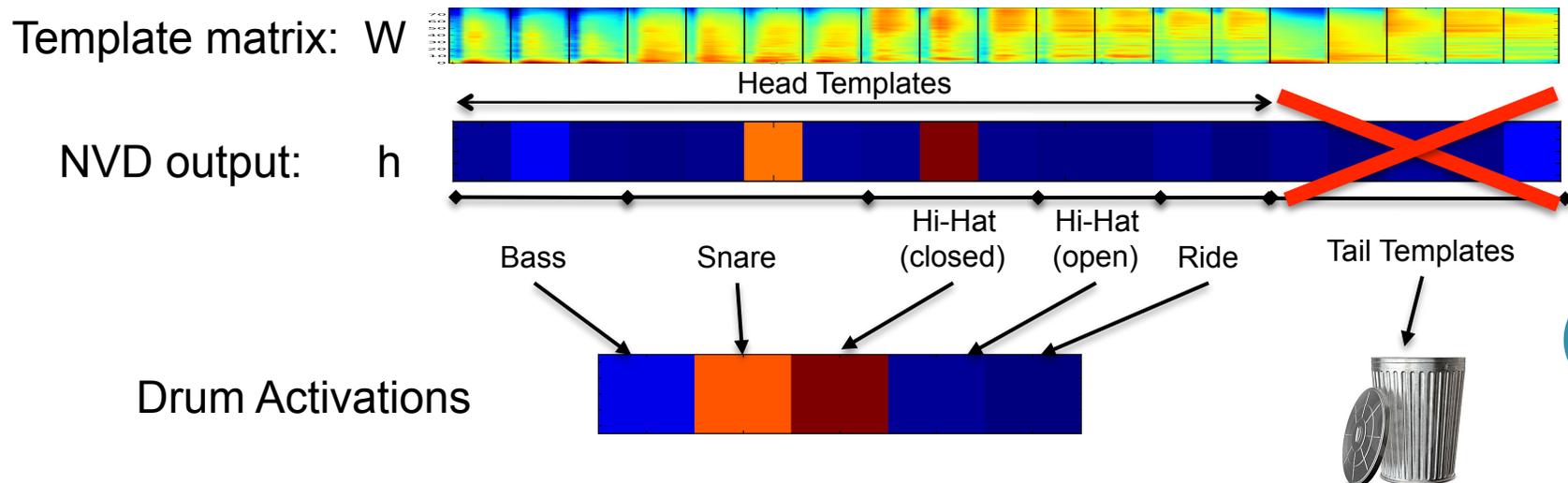
- Non-negative Vector Decomposition (NVD)
 - A simplification of Non-negative Matrix Factorization (NMF)
 - W matrix contains drum templates in its columns.
 - Adding a sparsity penalty (L1) on h improves NVD.

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



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.



EVALUATION

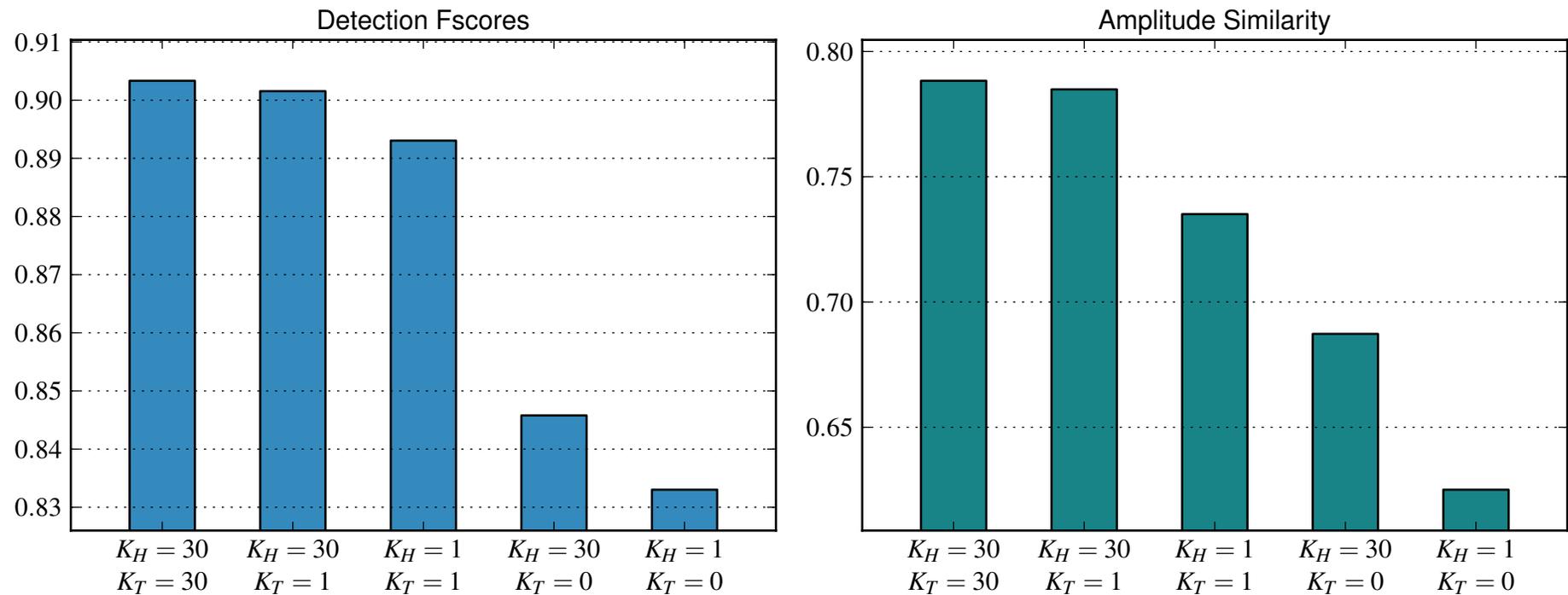
- Test Data:
 - 23 minutes total, 8022 drum onsets
 - 8 different drums/cymbals:
 - Bass, snare, hi-hat (open/closed), ride, 2 toms, crash.
 - Recorded to stereo using multiple microphones.
 - 50-100 training hits per drum.
- Parameters to vary for testing:
 - Maximum number of templates per drum {0,1,30}
- Result Metrics
 - Detection accuracy:
 - F-Score
 - Amplitude Fidelity
 - Cosine similarity

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

DETECTION RESULTS

- Varying maximum templates per drum.
- Adding tail templates helps the most.
- > 1 head template helps too.

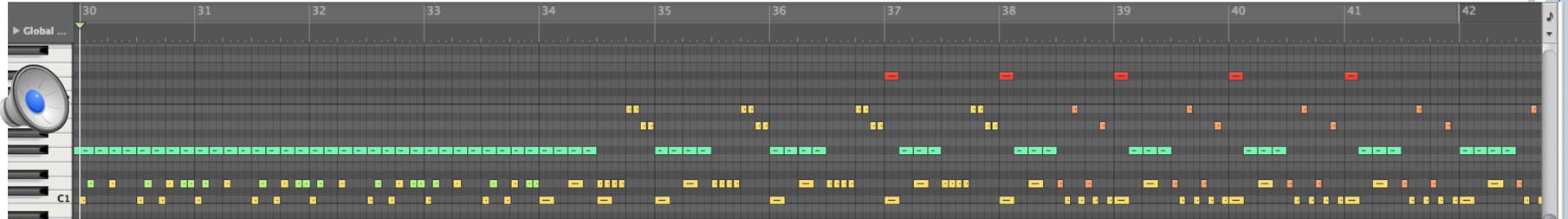
Detection Results Varying Max Templates per Drum



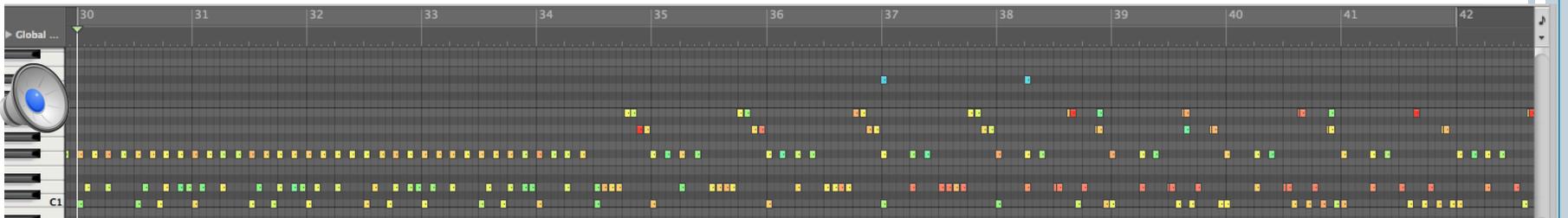
AUDIO EXAMPLES

- 100 BPM, rock, syncopated snare drum, fills/flams.
- Note the messy fills in the KH=1, KT=0 version

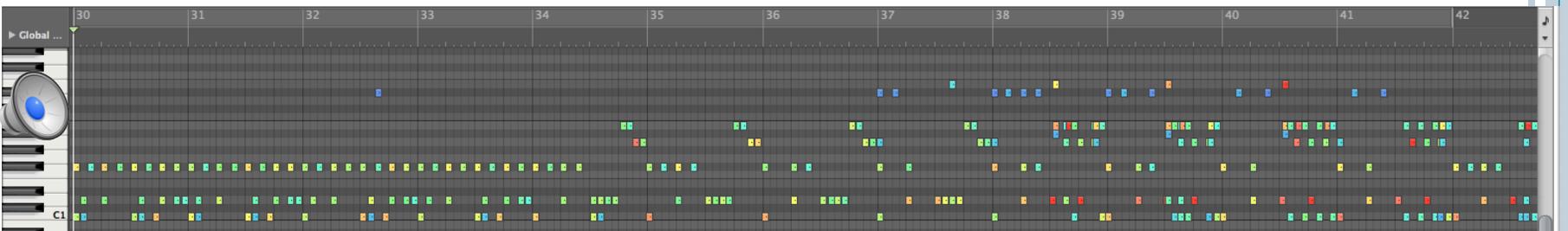
Original Performance



KH=30, KT=30



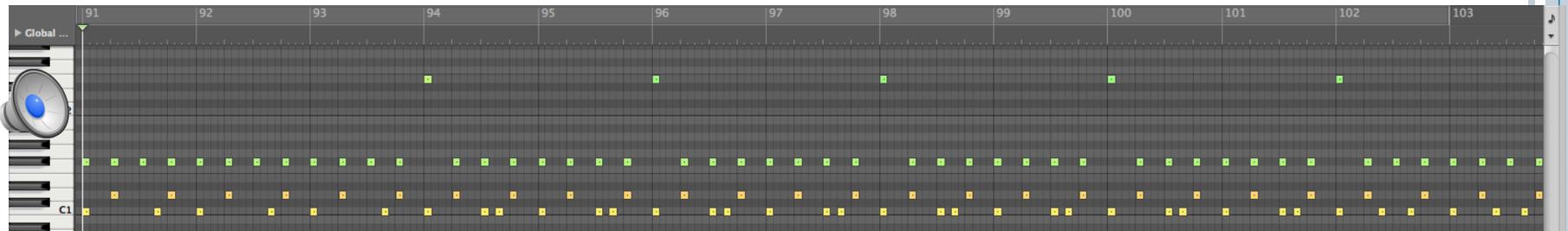
KH=1, KT=0



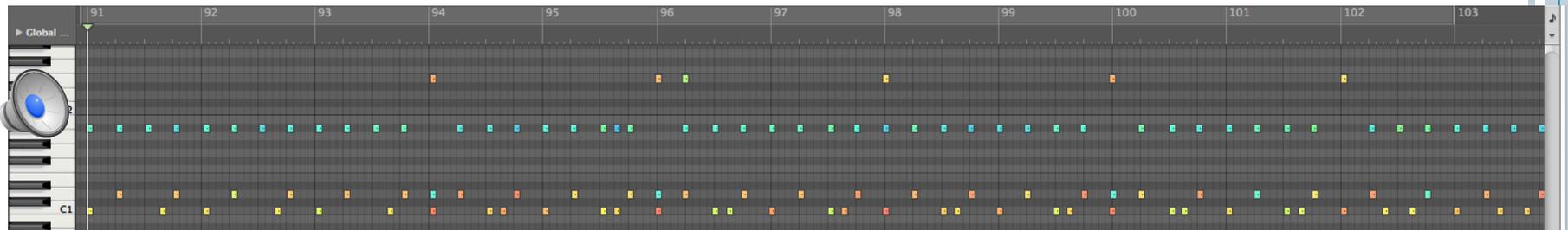
AUDIO EXAMPLES

- 181 BPM, fast rock, open hi-hat
- Note the extra hi-hat notes in the $KH=1, KT=0$ version.

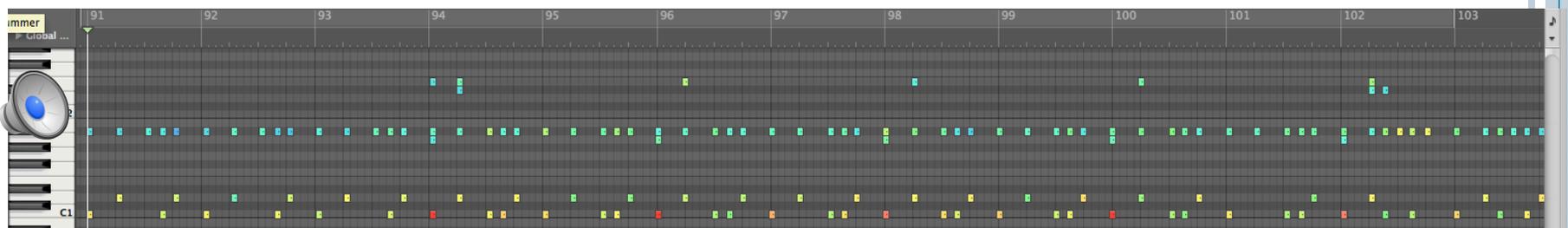
Original Performance



$KH=30, KT=30$



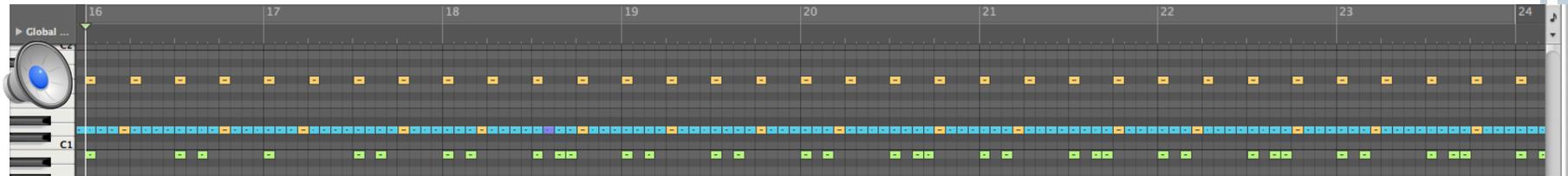
$KH=1, KT=0$



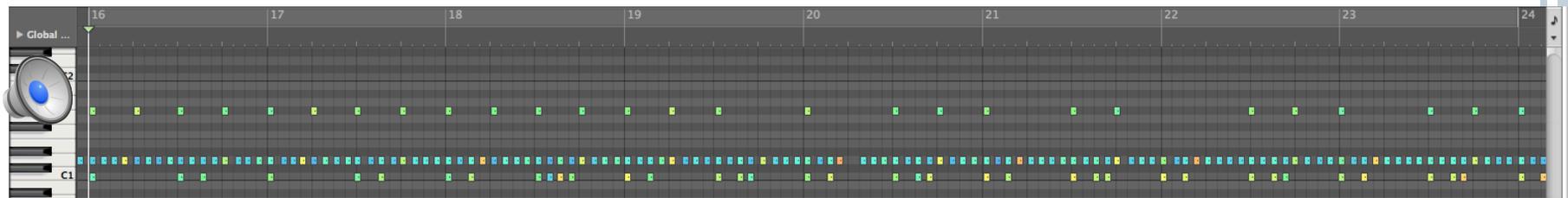
AUDIO EXAMPLES

- 94 BPM, snare drum march, accented notes
- Note the extra bass drum notes and many extra cymbals in the KH=1, KT=0 version.

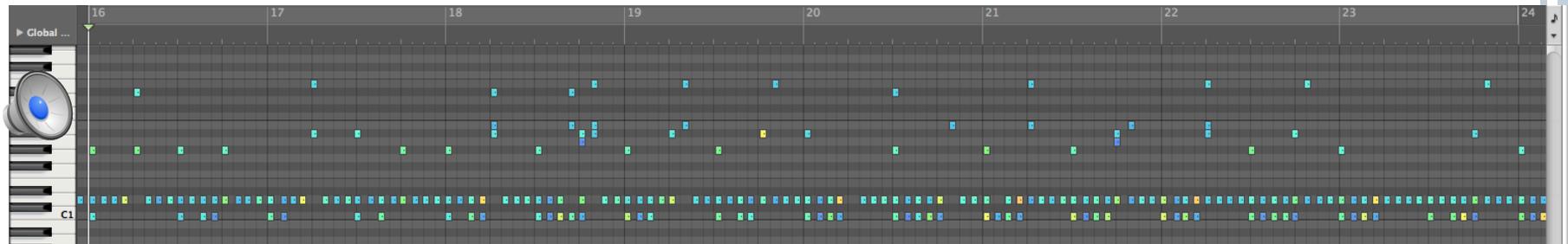
Original Performance



KH=30, KT=30



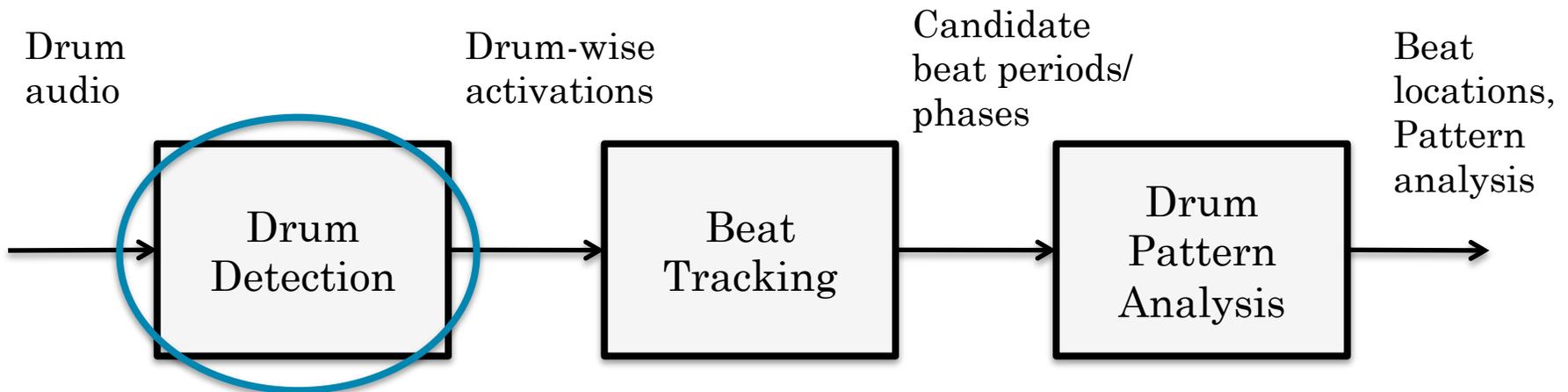
KH=1, KT=0



DRUM DETECTION SUMMARY

- Drum detection front end for a complete drum understanding system.
- Gamma Mixture Model
 - Cheaper to train than GMM (no covariance matrix)
 - More stable than GMM (no covariance)
 - Allows soft clustering with perceptual Itakura-Saito distance in the linear domain (important for learning NVD templates).
- Non-negative Vector Decomposition
 - Greatly improved with tail templates and multiple head templates per drum.
- Next steps
 - Online training of templates.
 - Training a more general model using many templates

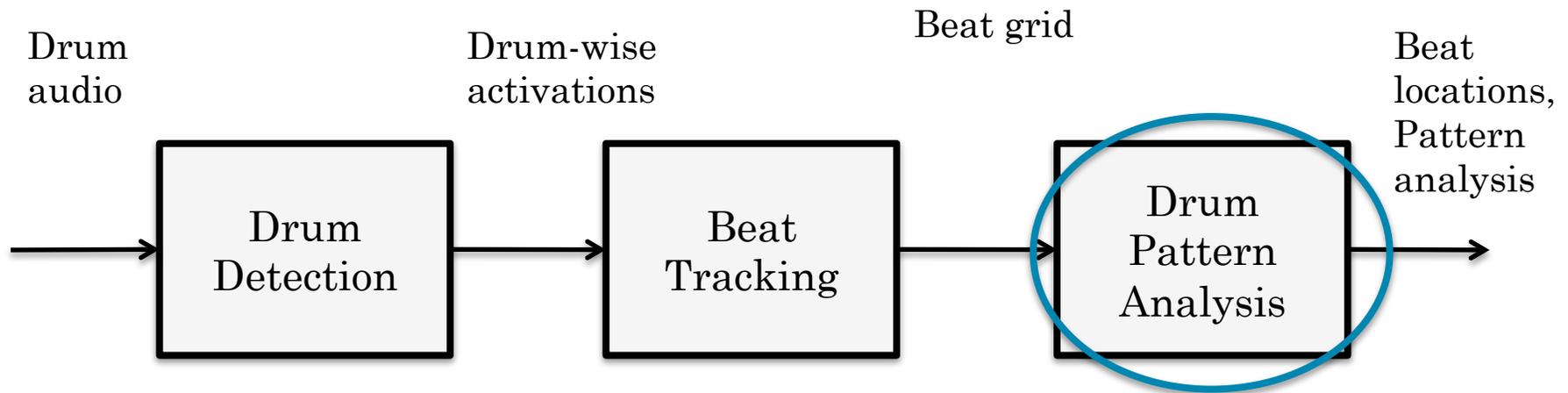
LIVE DRUM UNDERSTANDING SYSTEM



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

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

LIVE DRUM UNDERSTANDING SYSTEM

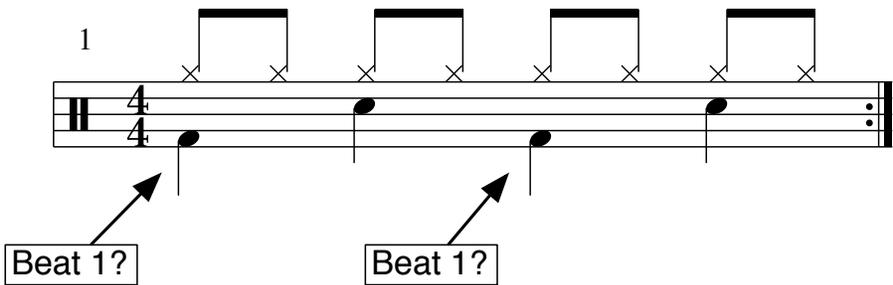


- Gamma Mixture Model training of drum templates
- Non-negative decomposition onto templates.

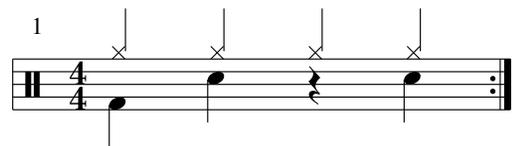
- Generative deep learning of drum patterns
- Stacked Conditional Restricted Boltzmann Machines

DRUM PATTERN ANALYSIS

- Desired information:
 - What style is this?
 - What is the meter? (4/4, 3/4...)
 - Double/half time feel?
 - **Where is the “one”?**

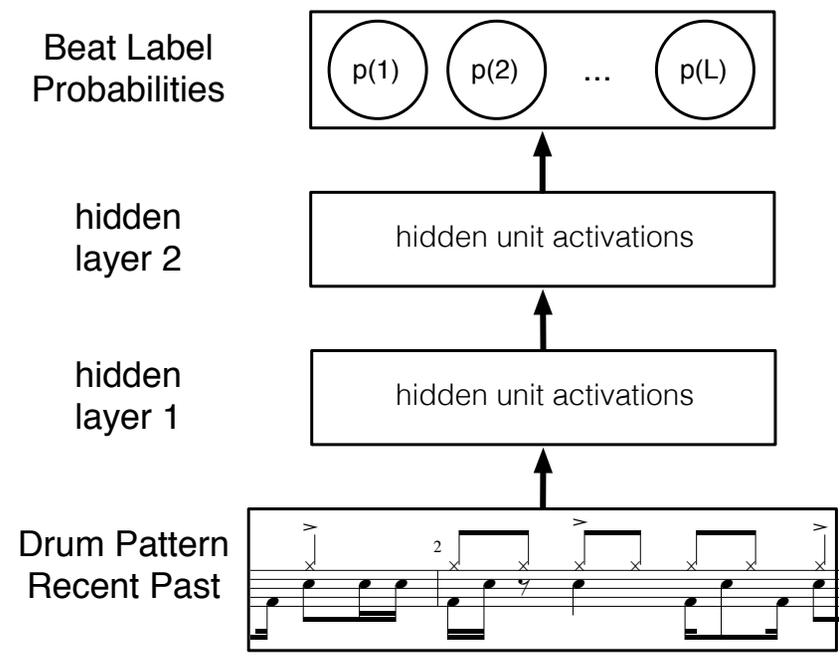


- Typical approach:
 - Drum pattern template correlation.
 - Align one or more templates with drum onsets.



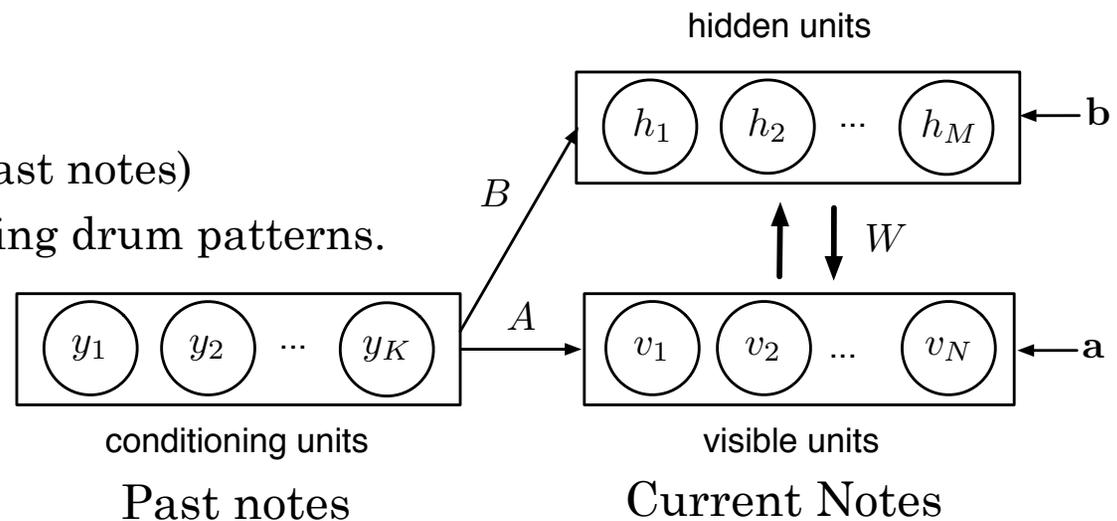
Instead...

Deep Neural Network



CONDITIONAL RESTRICTED BOLTZMANN MACHINE

- Modeling motion, language models, other sequences.
- Models: $P(v | y)$
- For music:
 - $P(\text{current notes} | \text{past notes})$
 - Useful for generating drum patterns.



- Used to pre-train neural network.
- Intuition for drums
 - Hidden units model the drummer's options given the recent past.
 - And therefore, code information about the state of the drumming.

TEST SETUP

- 173 twelve-measure sequences
 - Recorded on Roland V-Drums
 - 33,216 beat subdivisions (16 per measure).
 - Rock, funk, drum 'n' bass, metal, Brazilian rhythms.
- Network configurations
 - Each with 2 measures of context
 - ~90,000 weights each
 - Conditional RBM (CRBM) is a variation on the RBM
- 1. CRBM (3 layers)
 - 800 hidden units
- 2. CRBM→RBM (4 layers)
 - 600, 50 hidden units
- 3. CRBM→CRBM (4 layers)
 - 200, 50 hidden units
- 4. CRBM→CRBM→RBM (5 layers)
 - 200, 25, 25 hidden units

TRAINING

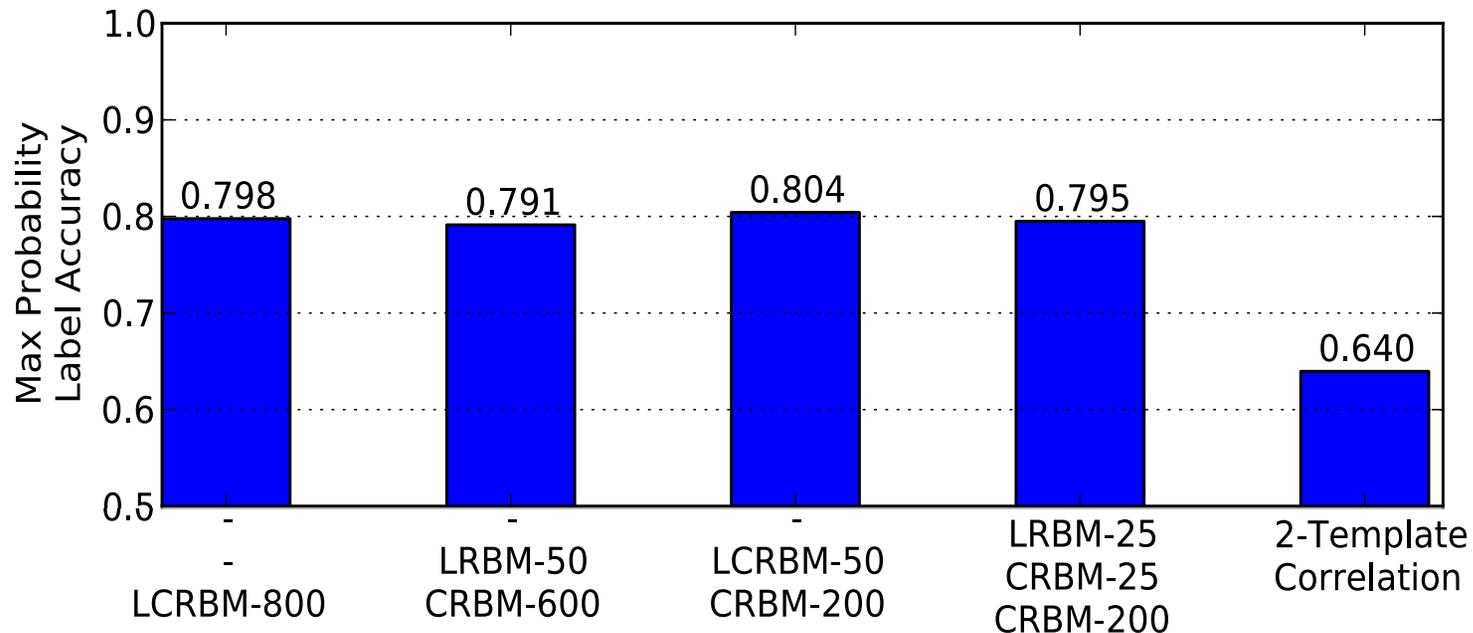
- Training tricks
 - L2 weight decay, momentum, dynamic learning rate.
- Implementation
 - Python with Gnumpy (numerical routines on GPU)
 - GPU computing very important to contemporary neural network research.
 - Around 30 minutes to train one model.



MEASURE ALIGNMENT RESULTS

- 3-fold cross-validated results.
- Neural network models eliminate half the errors of template correlation.
- Not a significant difference between NN models.
- This will change with a larger, more diverse dataset.

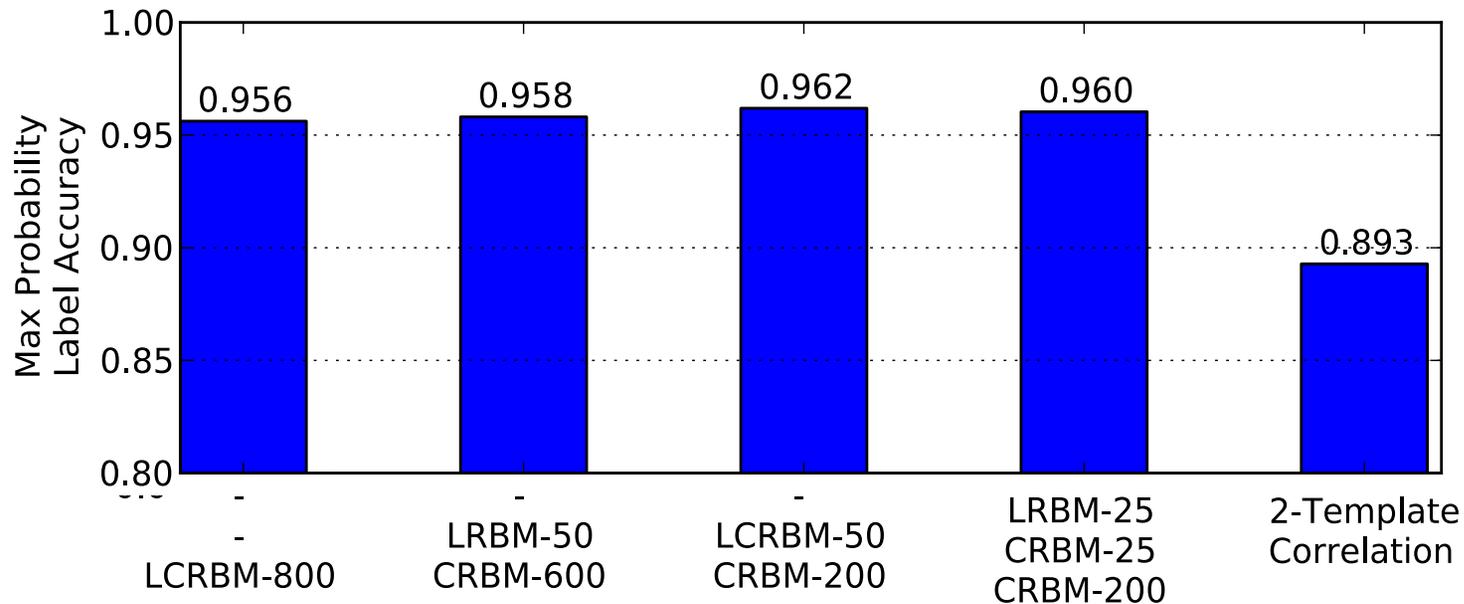
Measure (Whole Note) Alignment Accuracy by Model



QUARTER NOTE ALIGNMENT RESULTS

- Can compute quarter note alignment from full measure alignment probabilities.
- Quarter note alignment can help correct a beat tracker when it gets out of phase.
- Again, half the errors are eliminated.

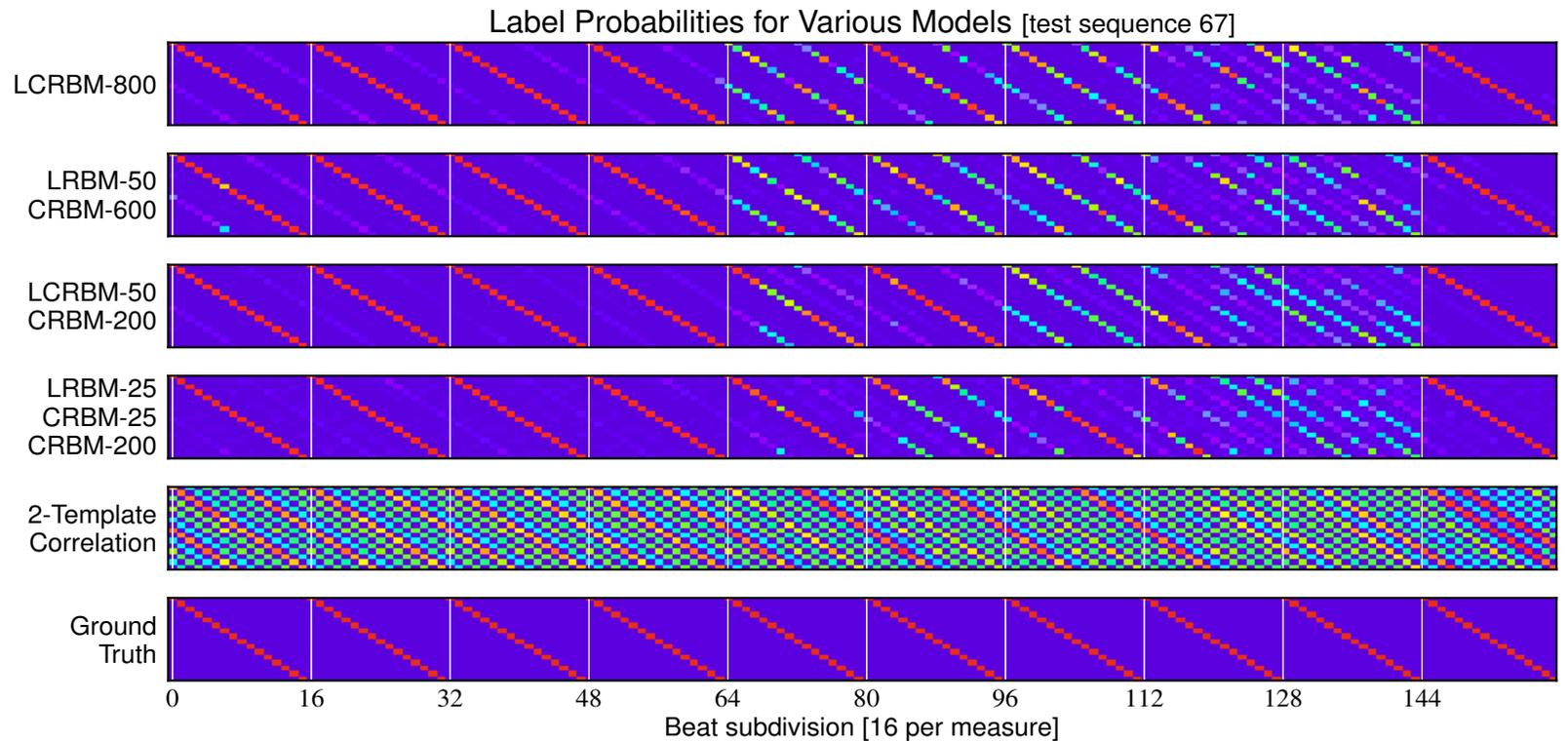
Quarter Note Alignment Accuracy by Model





EXAMPLE OUTPUT

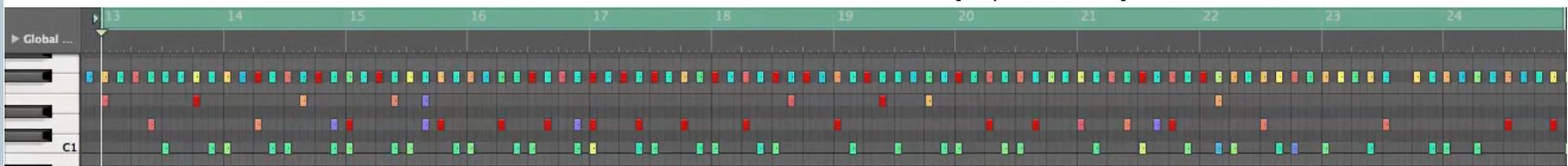
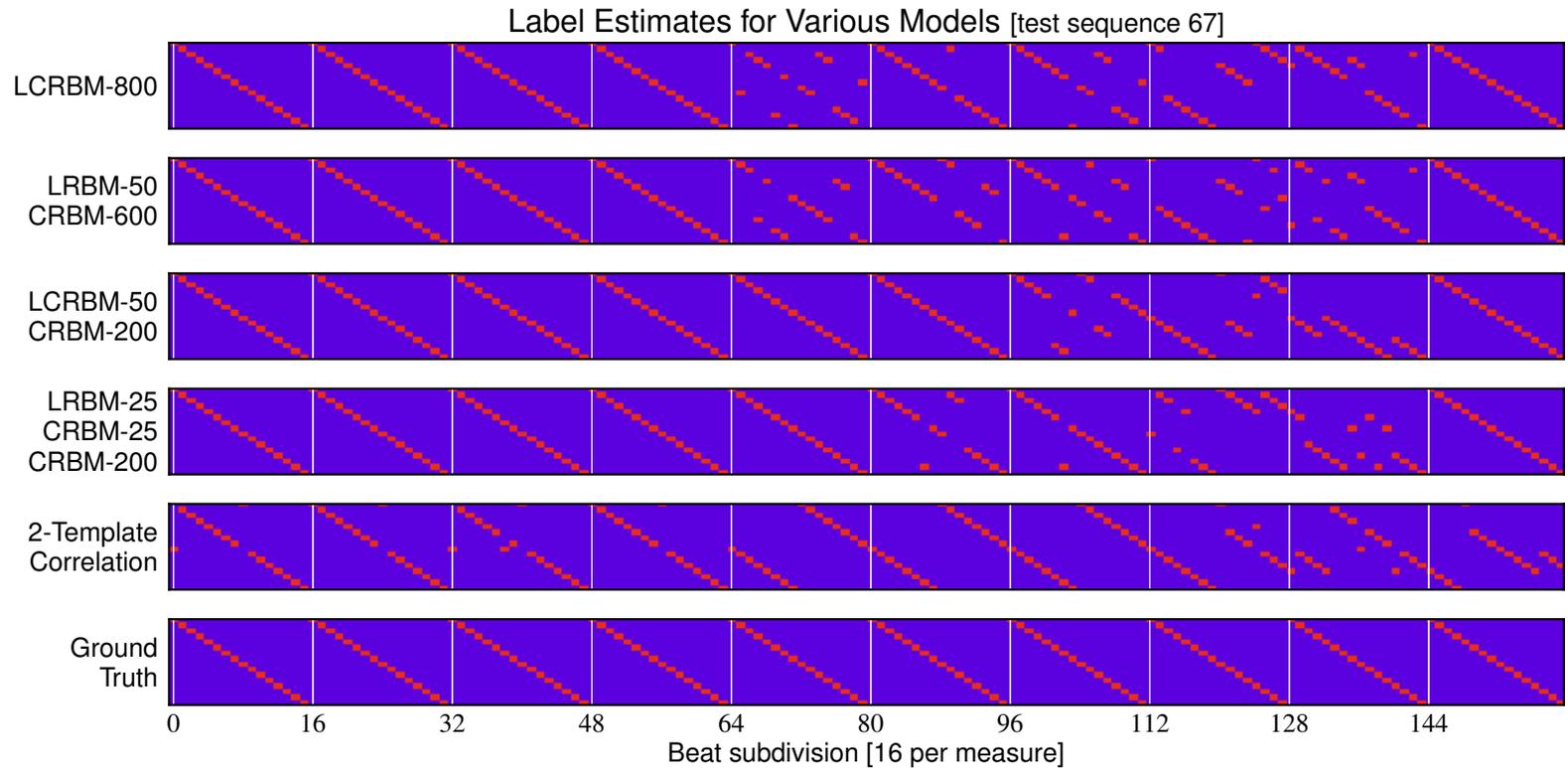
- Label Probabilities (Red = 1, Blue = 0)
- White lines denote measure boundaries
- Each column shows the probability of each label at that time.





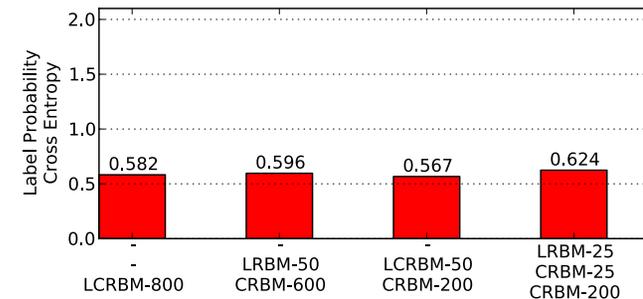
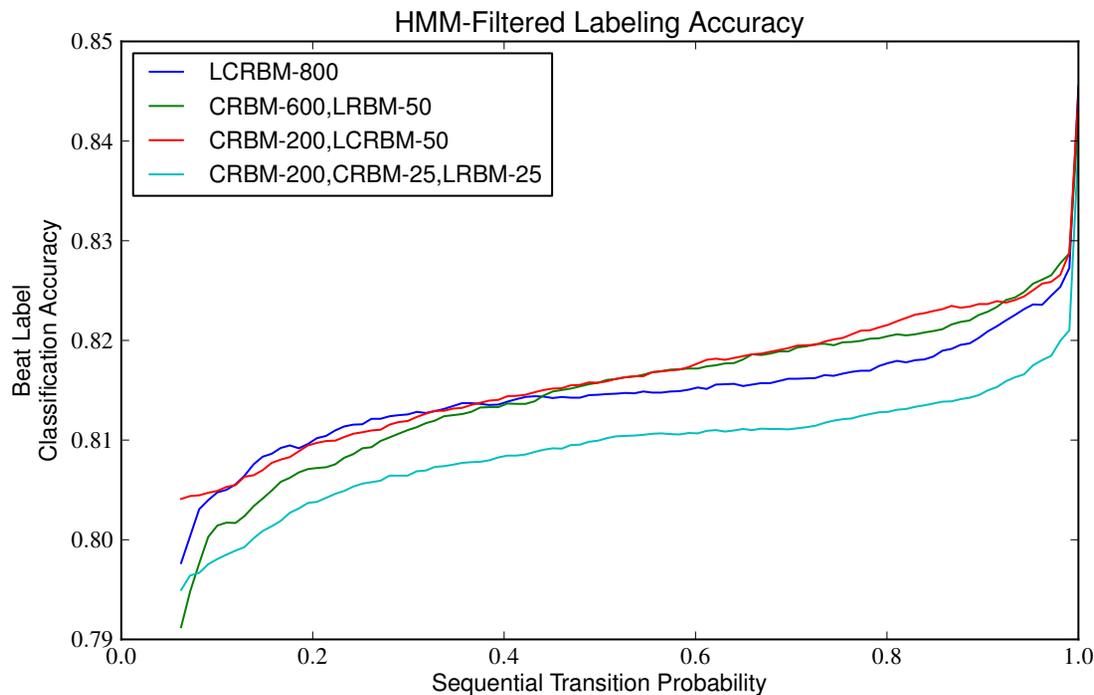
EXAMPLE OUTPUT

- Label Estimates (Red denotes estimate)
- White lines denote measure boundaries



HMM FILTERING

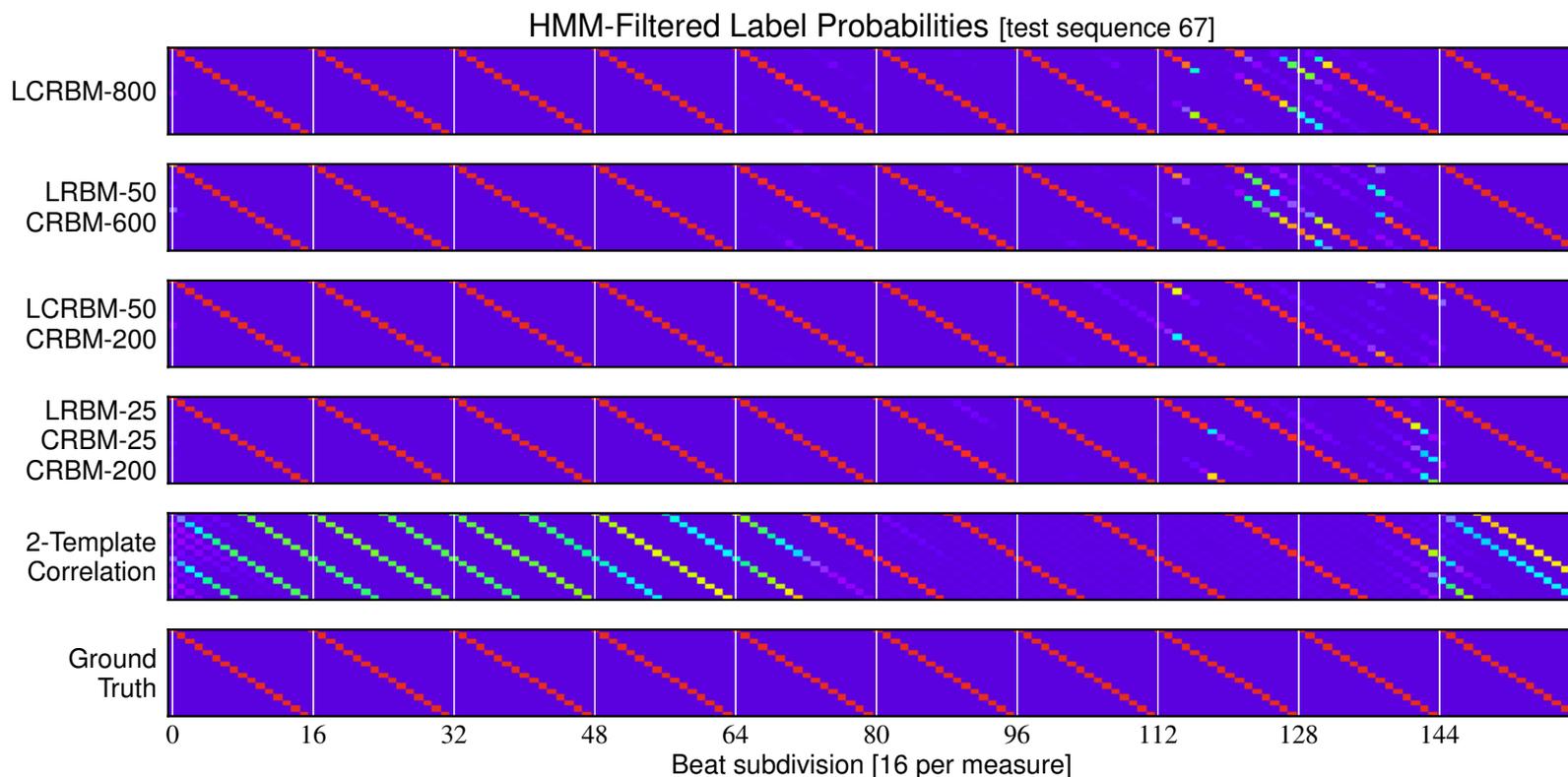
- Can use label probabilities as observation posteriors in HMM.
- Assign high probability to sequential label transitions.
- Models producing lower cross-entropy improve more when using HMM-filtering.





EXAMPLE OUTPUT

- HMM-Filtered Label Probabilities



ANALYSIS OF RESULTS

- Generatively pre-trained neural network models eliminate about half the errors compared to a baseline template correlation method.
- HMM-filtering can be used to improve accuracy.
- Overfitting present in backpropagation.
- Address overfitting with:
 - Larger dataset
 - “Dropout” [Hinton, 2013]
 - Many other regularization techniques
- Next steps:
 - Important next step is evaluation with much larger, more diverse dataset.
 - Evaluate ability to do style, meter classification.

SUMMARY

- Content-Based Music Information Retrieval
 - Mood, genre, onsets, beats, transcription, recommendation, and much, much more!
 - Exciting directions include feature learning and RNNs.
- Drum Detection
 - Learn multiple templates per drum using agglomerative gamma mixture model training.
 - The use of “tail” templates reduce false positives.
- Drum Pattern Analysis
 - A deep neural network can be used to make all kinds of rhythmic inferences.
 - For measure alignment, reduces errors by ~50% compared to template correlation.
 - Generalization can be improved using more data and regularization techniques during backpropagation.

GETTING INVOLVED IN MUSIC INFORMATION RETRIEVAL

- Check out the proceedings of ISMIR (free online):
 - <http://www.ismir.net/>
- Participate in MIREX (annual MIR eval):
 - http://www.music-ir.org/mirex/wiki/MIREX_HOME
- Join the Music-IR mailing list:
 - <http://listes.ircam.fr/wws/info/music-ir>
- Join the Music Information Retrieval Google Plus community (just started it):
 - <https://plus.google.com/communities/109771668656894350107>

REFERENCES

○ Genre/Mood Classification:

- [1] P. Hamel and D. Eck, “Learning features from music audio with deep belief networks,” Proc. of the 11th International Society for Music Information Retrieval Conference (ISMIR 2010), pp. 339–344, 2010.
- [2] M. Henaff, K. Jarrett, K. Kavukcuoglu, and Y. LeCun, “Unsupervised learning of sparse features for scalable audio classification,” Proceedings of International Symposium on Music Information Retrieval (ISMIR’11), 2011.
- [3] J. Anden and S. Mallat, “Deep Scattering Spectrum,” arXiv.org. 2013.
- [4] E. Schmidt and Y. Kim, “Learning rhythm and melody features with deep belief networks,” ISMIR, 2013.

REFERENCES

○ Onset Detection:

- [1] J. Bello, L. Daudet, S. A. Abdallah, C. Duxbury, M. Davies, and M. Sandler, “A tutorial on onset detection in music signals,” *IEEE Transactions on Speech and Audio Processing*, vol. 13, no. 5, p. 1035, 2005.
- [2] A. Klapuri, A. Eronen, and J. Astola, “Analysis of the meter of acoustic musical signals,” *IEEE Transactions on Speech and Audio Processing*, vol. 14, no. 1, p. 342, 2006.
- [3] J. Bello, C. Duxbury, M. Davies, and M. Sandler, “On the use of phase and energy for musical onset detection in the complex domain,” *IEEE Signal Processing Letters*, vol. 11, no. 6, pp. 553–556, 2004.
- [4] S. Böck, A. Arzt, F. Krebs, and M. Schedl, “Online realtime onset detection with recurrent neural networks,” presented at the International Conference on Digital Audio Effects (DAFx-12), 2012.
- [5] F. Eyben, S. Böck, B. Schuller, and A. Graves, “Universal onset detection with bidirectional long-short term memory neural networks,” *Proc. of ISMIR*, 2010.

REFERENCES

- Neural Networks :
 - [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet classification with deep convolutional neural networks,” *Advances in neural information processing systems*, 2012.
- Unsupervised Feature Learning:
 - [2] G. E. Hinton and S. Osindero, “A fast learning algorithm for deep belief nets,” *Neural Computation*, 2006.
 - [3] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, “Extracting and composing robust features with denoising autoencoders,” pp. 1096–1103, 2008.
- Recurrent Neural Networks:
 - [4] I. Sutskever, “Training Recurrent Neural Networks,” 2013.
 - [6] D. Eck and J. Schmidhuber, “Finding temporal structure in music: Blues improvisation with LSTM recurrent networks,” pp. 747–756, 2002.
 - [7] S. Bock and M. Schedl, “Polyphonic piano note transcription with recurrent neural networks,” pp. 121–124, 2012.
 - [8] N. Boulanger-Lewandowski, Y. Bengio, and P. Vincent, “Modeling temporal dependencies in high-dimensional sequences: Application to polyphonic music generation and transcription,” *arXiv preprint arXiv:1206.6392*, 2012.
 - [9] G. Taylor and G. E. Hinton, “Two Distributed-State Models For Generating High-Dimensional Time Series,” *Journal of Machine Learning Research*, 2011.

REFERENCES

- Drum Understanding:

- [1] E. Battenberg, “Techniques for Machine Understanding of Live Drum Performances,” PhD Thesis, University of California, Berkeley, 2012.
- [2] E. Battenberg and D. Wessel, “Analyzing drum patterns using conditional deep belief networks,” presented at the International Society for Music Information Retrieval Conference, Porto, Portugal, 2012.
- [3] E. Battenberg, V. Huang, and D. Wessel, “Toward live drum separation using probabilistic spectral clustering based on the Itakura-Saito divergence,” presented at the Audio Engineering Society Conference: 45th International Conference: Applications of Time-Frequency Processing in Audio, 2012, vol. 3.

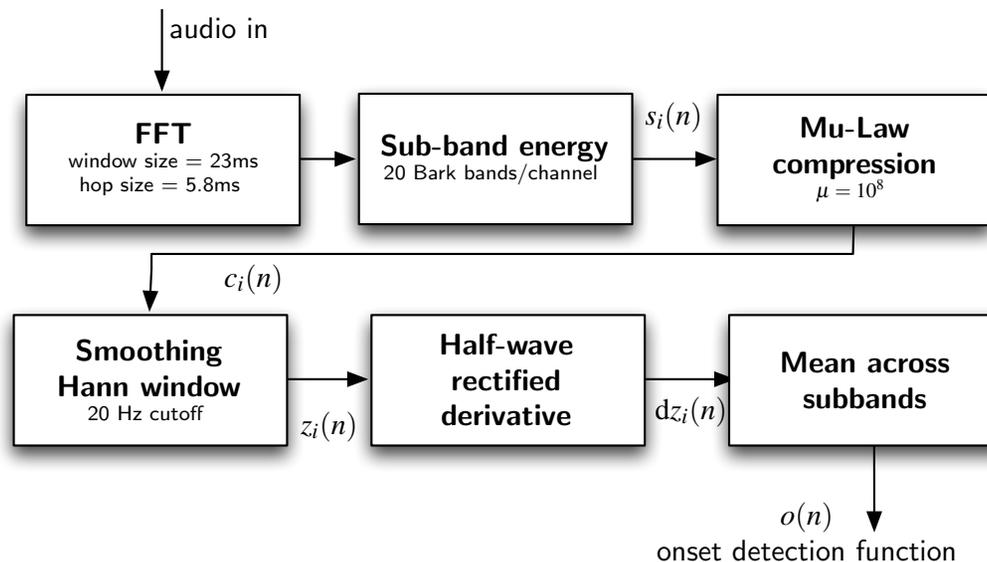
THANK YOU!



EXTRA SLIDES

ONSET DETECTION

- Onset Detection Function (ODF): Differentiated log-energy of multiple perceptual sub-bands.



- Onsets are located using ODF and a dynamic, causal peak-picking threshold.

DRUM SEPARATION

- Some Approaches to “Drum Transcription”
 - Feature-based classification [Gouyon 2001]
 - NMF with general drum templates [Paulus 2005]
 - General “average” drum templates.
 - Match and Adapt [Yoshii 2007]
 - Offline, iterative
- Requirements for “Drum Separation”
 - Online/Live operation
 - Useful with any percussion setup.
 - Which drum is playing when?
And at what dynamic level?



GAMMA DISTRIBUTION

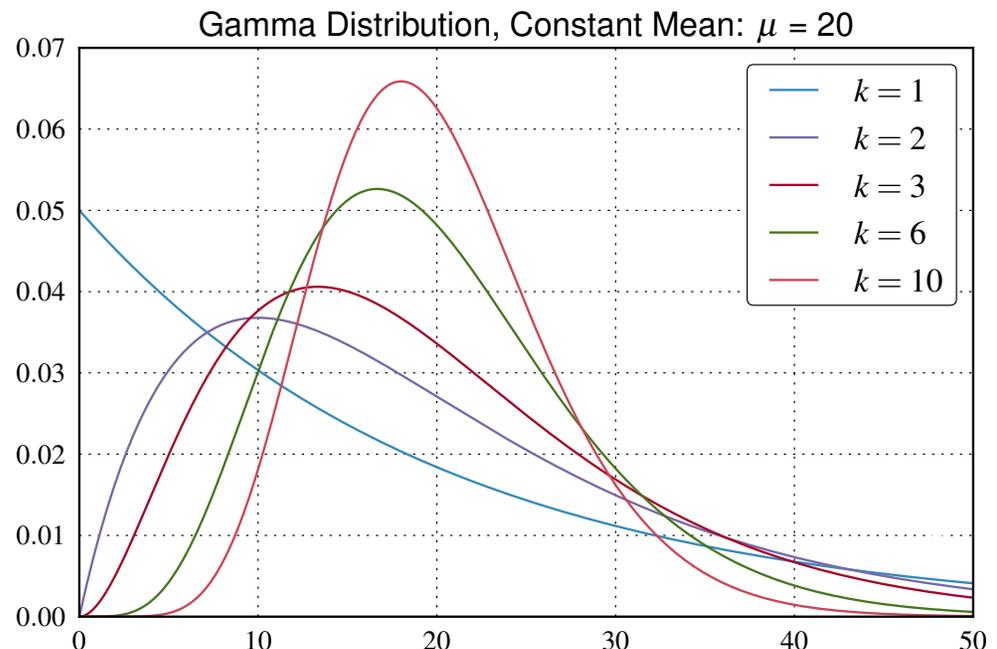
- Mixture model is composed of gamma distributions.
- The gamma distribution models the sum of k independent exponential distributions.

$$p(y|\lambda, k) = y^{k-1} \frac{\lambda^k e^{-\lambda y}}{\Gamma(k)},$$

$$E[y] = \mu = k/\lambda$$

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

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



GAMMA 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:

$$p(\vec{y}_n|\theta) = \sum_{l=1}^K \pi_l p(\vec{y}_n|\vec{\lambda}_l, k) \quad \theta = \{\vec{\lambda}_l, \pi_l\}_{l=1}^K$$
$$\pi_l = p(x_n = l)$$

THE EM ALGORITHM: GAMMA EDITION

- E-step: (compute posteriors)

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

- M-step: (update parameters)

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

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

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

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.

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

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

DECOMPOSING ONSETS ONTO TEMPLATES

- To solve this problem (add L1 penalty):

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

- I 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:

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

RESTRICTED BOLTZMANN MACHINE

- Stochastic autoencoder [Hinton 2006]

- Probabilistic graphical model
- Weights/biases define:

- $$p(\mathbf{v}, \mathbf{h}) = \frac{e^{-E(\mathbf{v}, \mathbf{h})}}{Z} \quad E(\mathbf{v}, \mathbf{h}) = -\mathbf{a}^T \mathbf{v} - \mathbf{b}^T \mathbf{h} - \mathbf{v}^T W \mathbf{h}$$

- Factorial conditionals:

$$p(\mathbf{v}|\mathbf{h}) = \prod_i p(v_i|\mathbf{h}) \quad p(\mathbf{h}|\mathbf{v}) = \prod_i p(h_i|\mathbf{v})$$

- Logistic activation probabilities

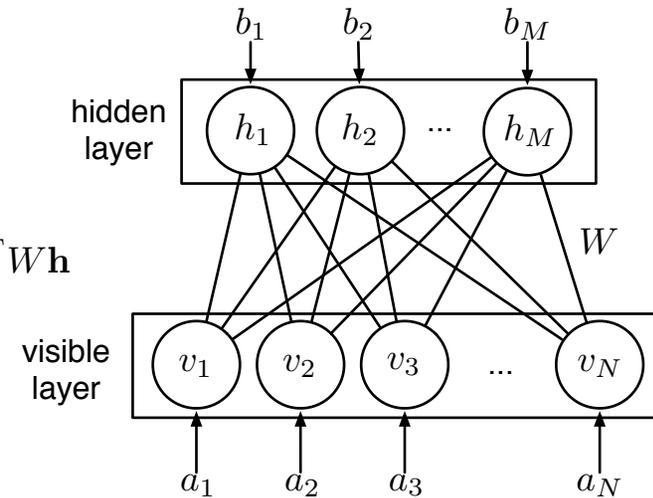
$$p(v_i = 1|\mathbf{h}) = \sigma(a_i + \sum_j W_{ij}h_j)$$

$$p(h_i = 1|\mathbf{v}) = \sigma(b_i + \sum_j W_{ji}v_j)$$

- Binary units (as opposed to real-valued units)

- Act as a strong regularizer (prevent overfitting)

- Training: maximize likelihood of data under the model $p(\mathbf{v})$



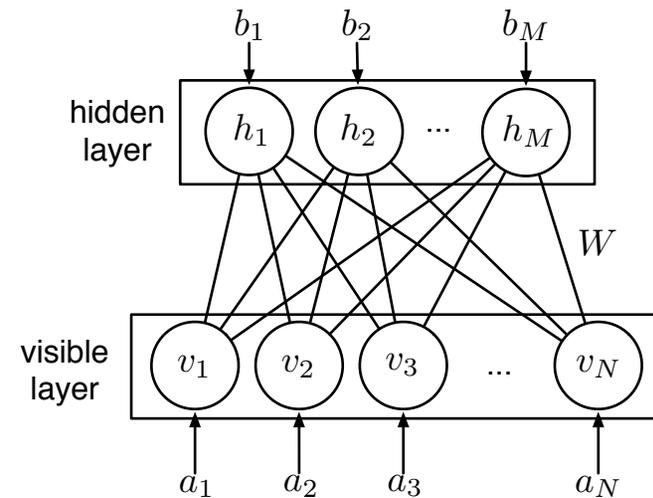
CONTRASTIVE DIVERGENCE

- ML learning can be done using Gibbs sampling to compute samples from:
 - The joint $p(\mathbf{v}, \mathbf{h})$
 - By taking alternating samples from:
 $p(v_i = 1 | \mathbf{h}) = \sigma(a_i + \sum_j W_{ij} h_j)$
 $p(h_i = 1 | \mathbf{v}) = \sigma(b_i + \sum_j W_{ji} v_j)$
- But, this can take a very long time to converge.
- Approximation:
“Contrastive Divergence” [Hinton 2006]
 - Take only a few alternating samples (k) for each update. (CD- k)

$$\Delta W_{ij} \propto \langle v_i h_j \rangle_0 - \langle v_i h_j \rangle_k$$

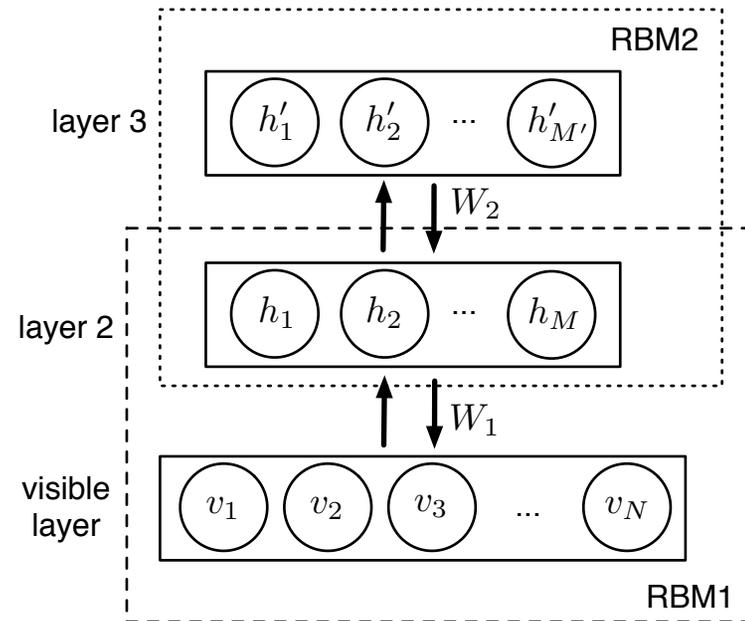
$$\Delta a_i \propto \langle v_i \rangle_0 - \langle v_i \rangle_k$$

$$\Delta b_j \propto \langle h_j \rangle_0 - \langle h_j \rangle_k$$



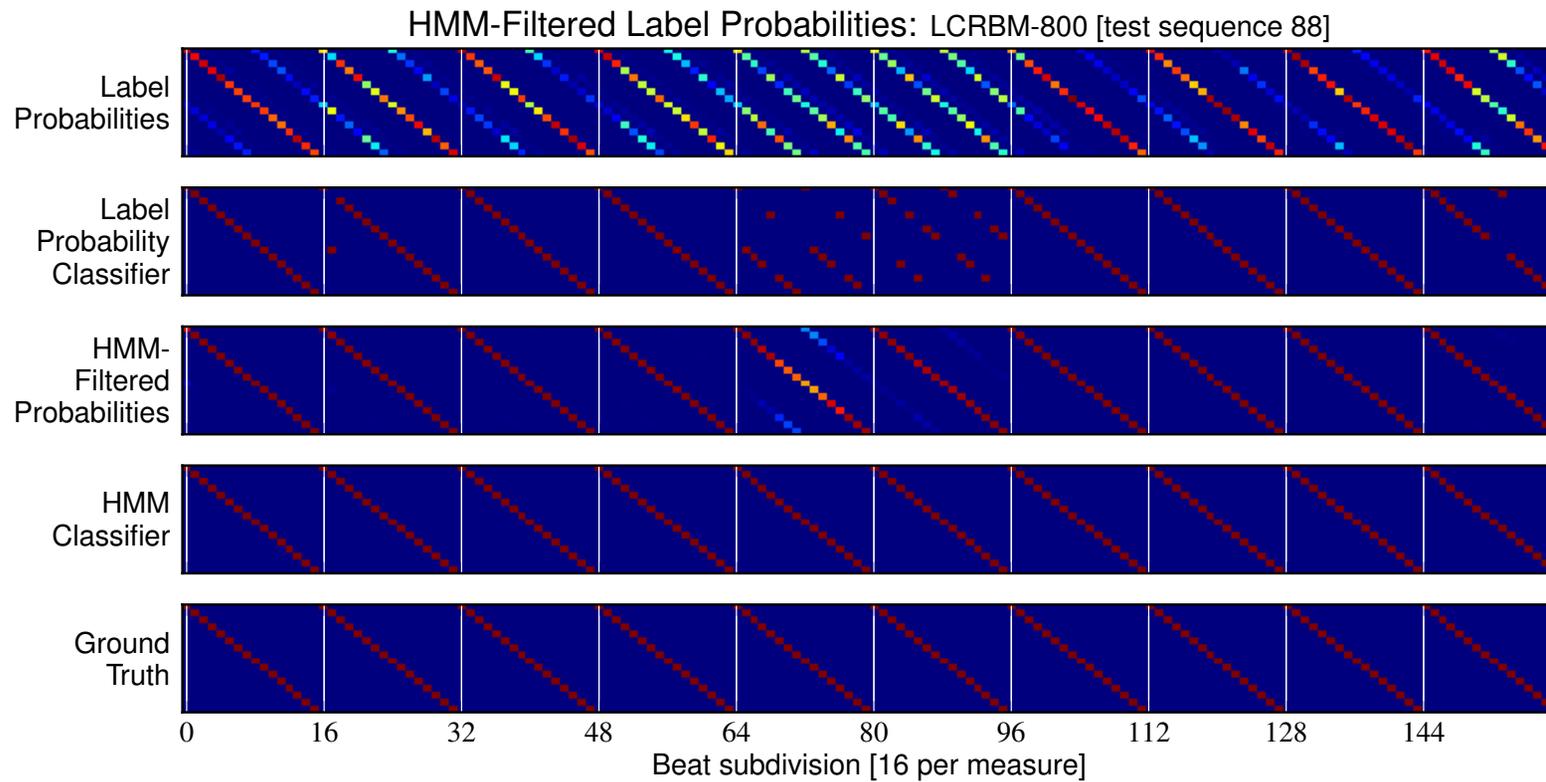
DEEP BELIEF NETWORKS (STACKING RBMs)

- Deep Belief Network (4+ layers)
 - Pre-train each layer as an RBM
- Greedy Pre-training
 - Train first level RBM
 - Use real-valued hidden unit activations of an RBM as input to subsequent RBM.
 - Train next RBM
 - Repeat until *deep* enough.
- Fine-Tuning
 - After pre-training, network is discriminatively fine-tuned using backpropagation of the cross-entropy error.



EXAMPLE HMM FILTERING

- White lines denote measure boundaries.





EXAMPLE OUTPUT

- HMM-Filtered Label Estimates

