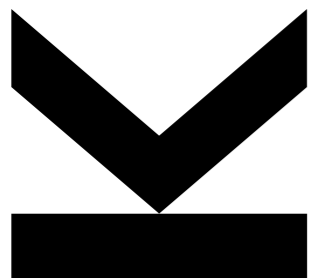


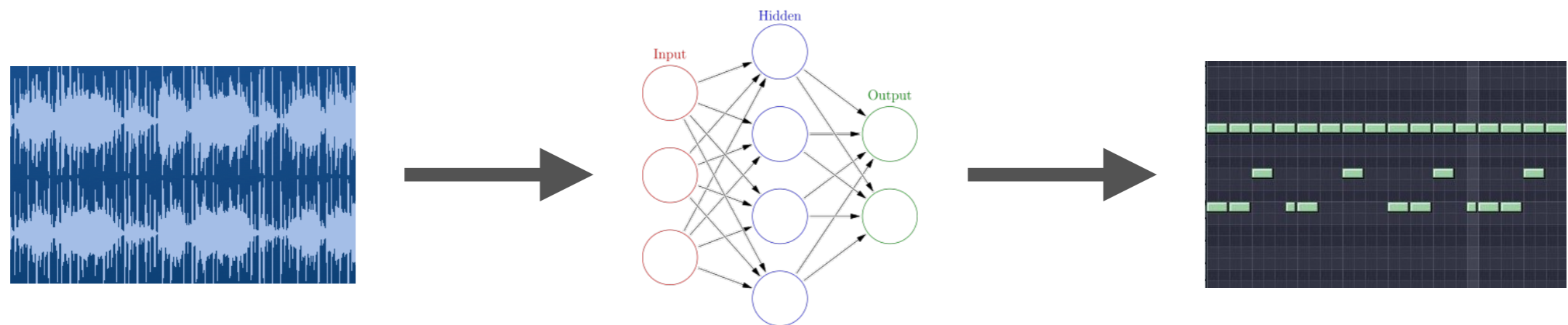
# DRUM TRANSCRIPTION FROM POLYPHONIC MUSIC WITH RECURRENT NEURAL NETWORKS



Richard Vogl<sup>1,2</sup>, Matthias Dorfer<sup>1</sup>, Peter Knees<sup>2</sup>

[richard.vogl@tuwien.ac.at](mailto:richard.vogl@tuwien.ac.at), [matthias.dorfer@jku.at](mailto:matthias.dorfer@jku.at), [peter.kness@tuwien.ac.at](mailto:peter.kness@tuwien.ac.at)

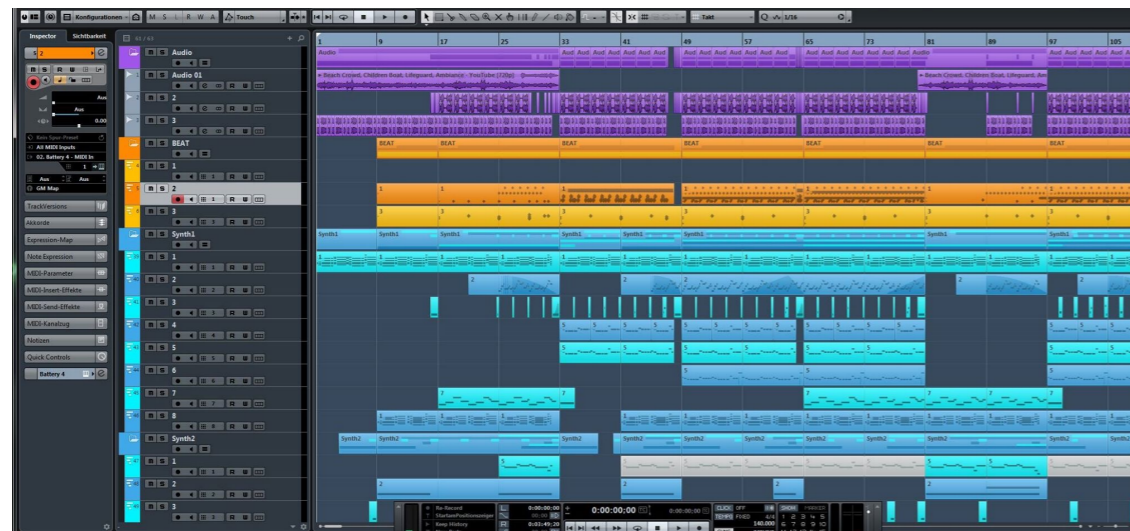
# INTRODUCTION



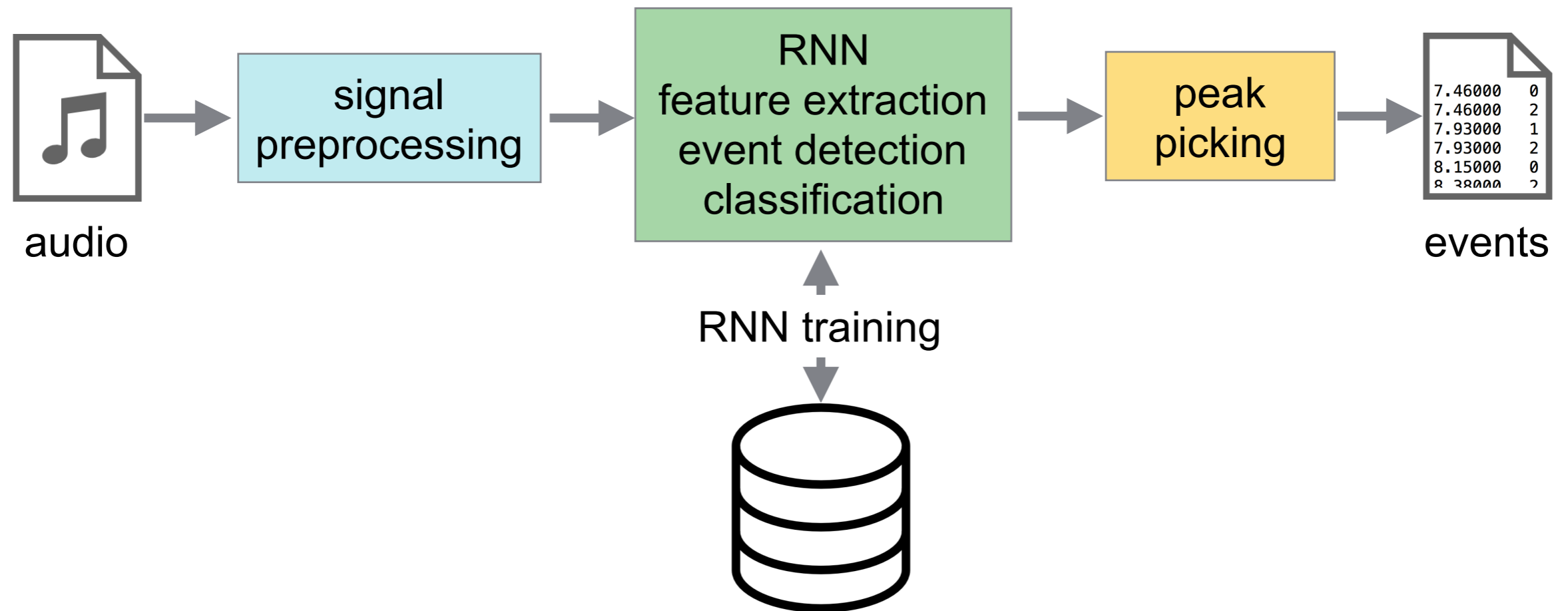
- Goal: **model** for drum note detection in polyphonic music
  - **In:** Western popular **music containing drums**
  - **Out:** **Symbolic representation** of notes played by drum instruments
- Focus on **three major drum instruments:** snare, bass drum, hi-hat

# INTRODUCTION

- Wide range of **applications**
  - **Sheet music** generation
  - **Re-synthesis** for music production
  - Higher level **MIR** tasks



# SYSTEM ARCHITECTURE



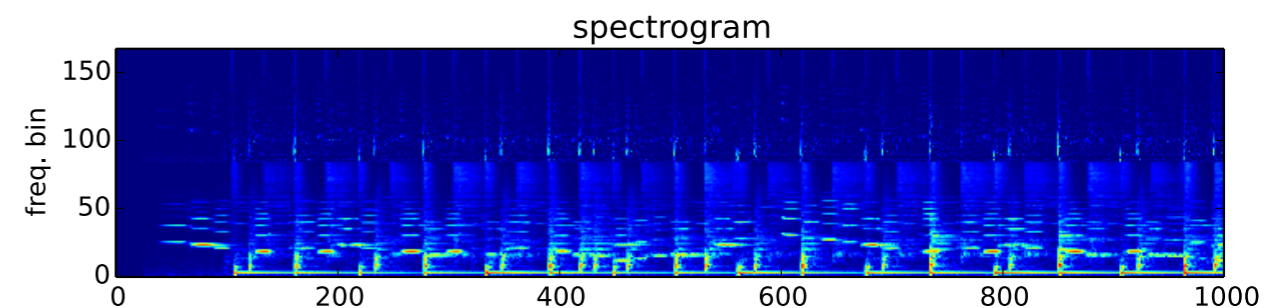
# ADVANTAGES OF RNNS

- Relatively **easy to fit** large and diverse datasets
- Once trained, **computational complexity** of transcription relatively **low**
- **Online capable**
- **Generalize well**
- **Easy to adapt** to new data
- **End-to-end**: learn features, event detection, and classification at once
- **Scale better** with number of instruments (rank problem in NMF)
- Trending topic: lots of **theoretical work** to benefit from

# DATA PREPARATION



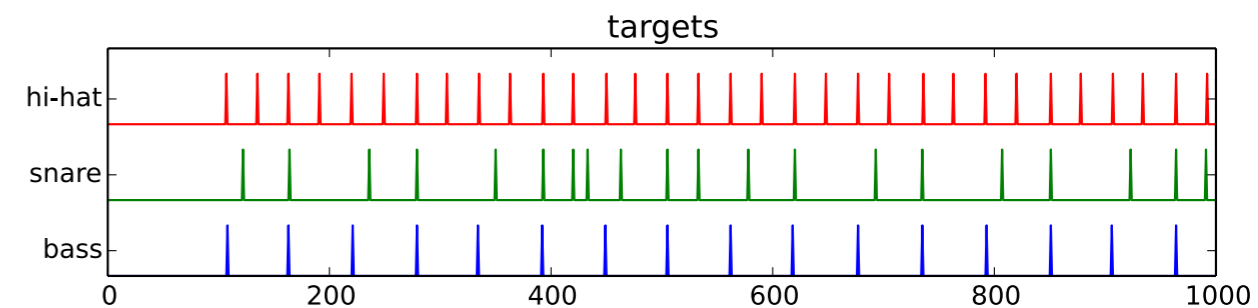
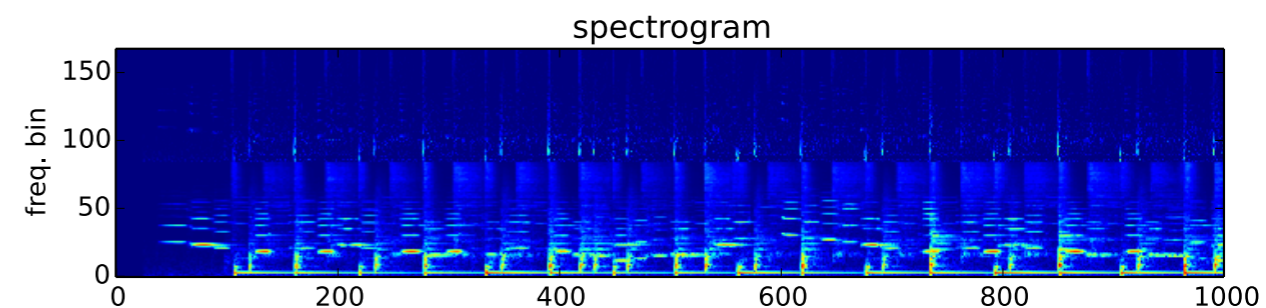
- Signal preprocessing
  - Log magnitude spectrogram @ 100Hz
  - Log frequency scale, 84 frequency bins
  - Additionally 1st order differential
  - 168 value input vector for RNN



# DATA PREPARATION



- Signal preprocessing
  - Log magnitude spectrogram @ 100Hz
  - Log frequency scale, 84 frequency bins
  - Additionally 1st order differential
  - 168 value input vector for RNN
- RNN targets
  - Annotations from training examples
  - Target vectors @ 100Hz frame rate

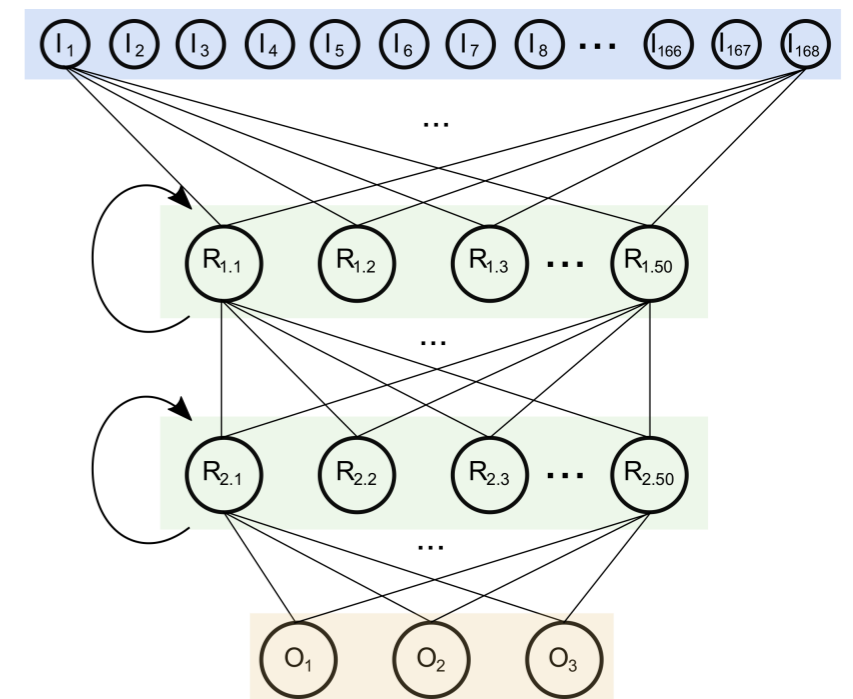




# RNN ARCHITECTURE

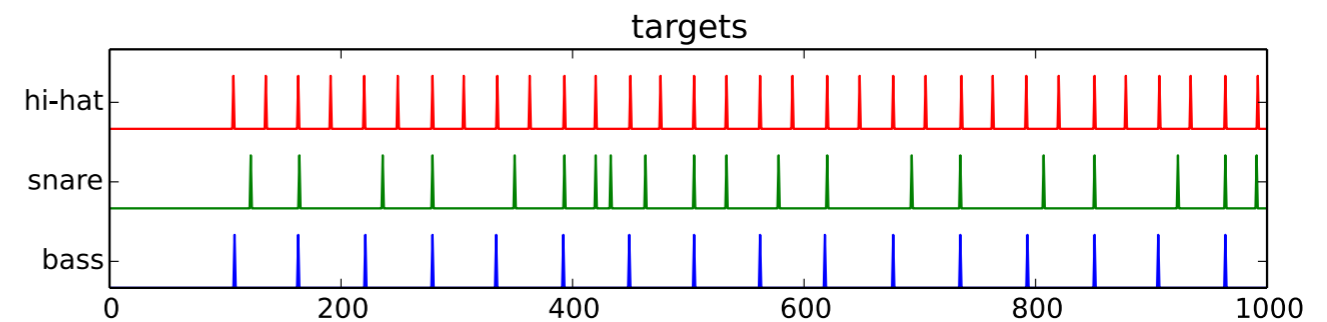
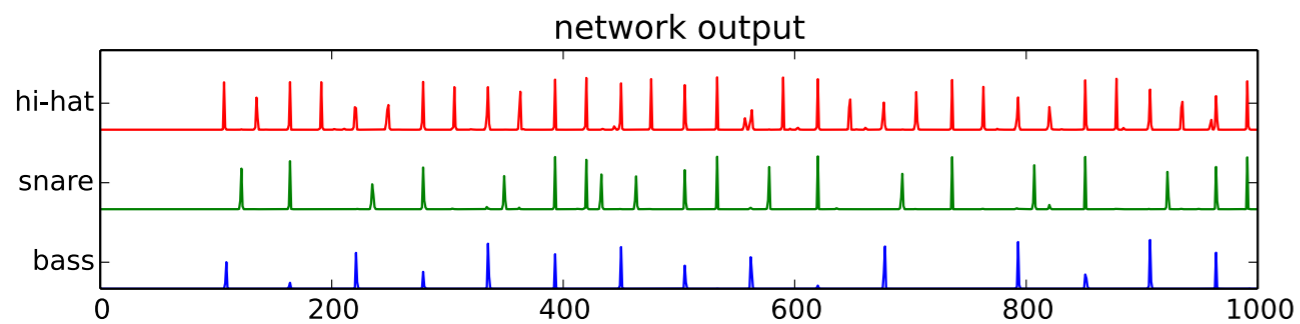


- Two layers containing 50 **GRUs** each
  - Recurrent connections
- Output: dense layer with three **sigmoid units**
  - No softmax: events are independent
  - Value represent certainty/pseudo-probability of drum onset
  - Does not model intensity/velocity





# PEAK PICKING



Select onsets at position  $n$  in activation function  $F(n)$  if:

$$F(n) = \max(F(n - m), \dots, F(n)),$$

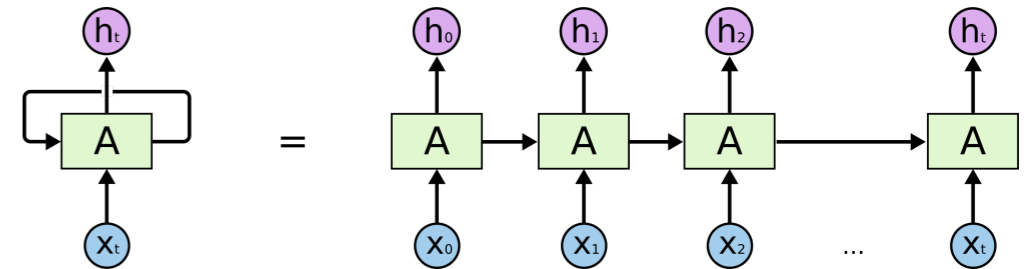
$$F(n) \geq \text{mean}(F(n - a), \dots, F(n)) + \delta,$$

$$n - n_{lp} > w,$$

[Böck et. al 2012]

# RNN TRAINING

- Backpropagation through time (**BPTT**)
- **Unfold RNN in time** for training
- Loss ( $\mathcal{L}$ ): mean **cross-entropy** between output ( $\hat{y}_n$ ) and targets ( $y_n$ ) for each instrument
- Mean over instruments with **different weighting** ( $w_i$ ) per instrument ( $\sim +3\%$  f-measure)
- Update model parameters ( $\theta$ ) using gradient ( $\mathcal{G}$ ) calculated on **mini-batch** and learn rate ( $\eta$ )



[Olah 2015]

$$\mathcal{L}(\Theta) = -\frac{1}{N} \sum_{n=1}^N [y_n \log \hat{y}_n + (1 - y_n) \log(1 - \hat{y}_n)]$$

$$\mathcal{L}(\Theta) = \sum_{i=1}^I w_i \mathcal{L}_i(\Theta) \quad \sum_{i=1}^I w_i = 1$$

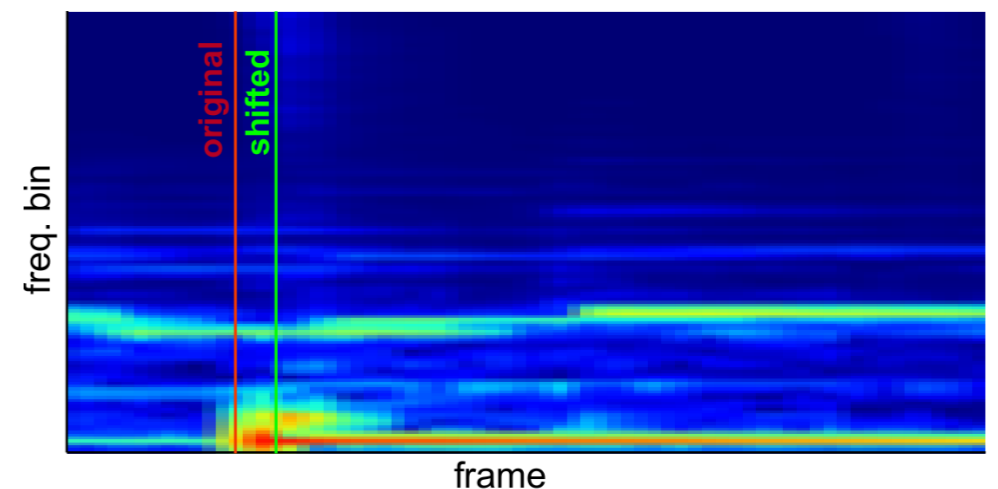
$$\mathcal{G}_t = \nabla_{\Theta} \mathcal{L}(\Theta_t)$$

$$\Theta_{t+1} = \Theta_t - \eta \mathcal{G}_t$$

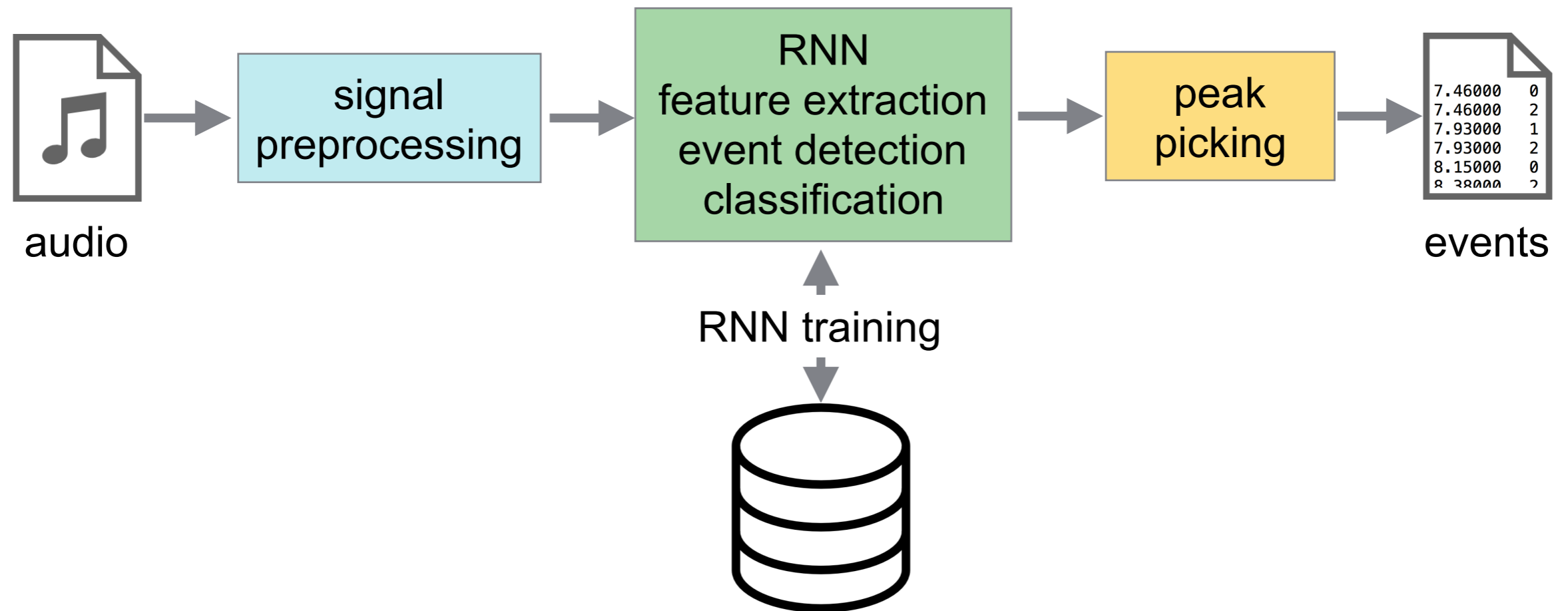
# RNN TRAINING (2)

- **RMSprop**
  - uses weight for learn rate based on moving mean squared gradient  $E[\mathcal{G}^2]$
- **Data augmentation**
  - Random transformations of training samples (pitch shift, time stretch)
- **Drop-out**
  - Randomly disable connections between second GRU layer and dense layer
- **Label time shift** instead of BDRNN

$$\cancel{\Theta_{t+1} = \Theta_t - \eta \mathcal{G}_t}$$
$$E[\mathcal{G}^2]_t = 0.9E[\mathcal{G}^2]_{t-1} + 0.1\mathcal{G}_t^2$$
$$\Theta_{t+1} = \Theta_t - \frac{\eta}{\sqrt{E[\mathcal{G}^2]_t + \epsilon}} \mathcal{G}_t$$



# SYSTEM ARCHITECTURE



# DATA / EVALUATION

- IDMT-SMT-Drums [Dittmar and Gärtner 2014]
  - Three classes (Real, Techno, and Wave / recorded/synthesized/sampled)
  - 95 simple solo drum tracks (30sec), plus training and single instrument tracks
- ENST-Drums [Gillet and Richard 2006]
  - Drum recordings, three drummers on three different drum kits
  - ~75 min per drummer, training, solo tracks plus accompaniment
- Precision, Recall, F-measure for drum note onsets
- Tolerance: 20ms

# EXPERIMENTS

- **SMT optimized**
  - Six fold cross-validation on randomized split of solo drum tracks
- **SMT solo**
  - Three fold cross-validation on different types of solo drum tracks
- **ENST solo**
  - Three fold cross-validation on solo drum tracks of different drummers / drum kits
- **ENST accompanied**
  - Three fold cross-validation on tracks with accompaniment

# RESULTS

Method	SMT opt.	SMT solo	ENST solo	ENST acc.
<b>NMF-SAB</b> [Dittmar and Gärtner 2014]	95.0	—	—	—
<b>PFNMF</b> [Wu and Lerch 2015]	—	81.6	77.9	72.2
<b>HMM</b> [Paulus and Klapuri 2009]	—	—	81.5	74.7
<b>BDRNN</b> [Southall et al. 2016]	96.1	83.3	73.2	66.9
<b><i>tsRNN</i></b>	<b>96.6</b>	<b>92.5</b>	<b>83.3</b>	<b>75.0</b>
	$\delta = 0.10$	$\delta = 0.15$	$\delta = 0.15$	$\delta = 0.10$



# RESULTS

Method	SMT opt.	SMT solo	ENST solo	ENST acc.
<b>NMF-SAB</b> [Dittmar and Gärtner 2014]	95.0	—	—	—
<b>PFNMF</b> [Wu and Lerch 2015]	—	81.6	77.9	72.2
<b>HMM</b> [Paulus and Klapuri 2009]	—	—	81.5	74.7
<b>BDRNN</b> [Southall et al. 2016]	96.1	83.3	73.2	66.9
<b><i>tsRNN</i></b>	<b>96.6</b>	<b>92.5</b>	<b>83.3</b>	<b>75.0</b>
	$\delta = 0.10$	$\delta = 0.15$	$\delta = 0.15$	$\delta = 0.10$

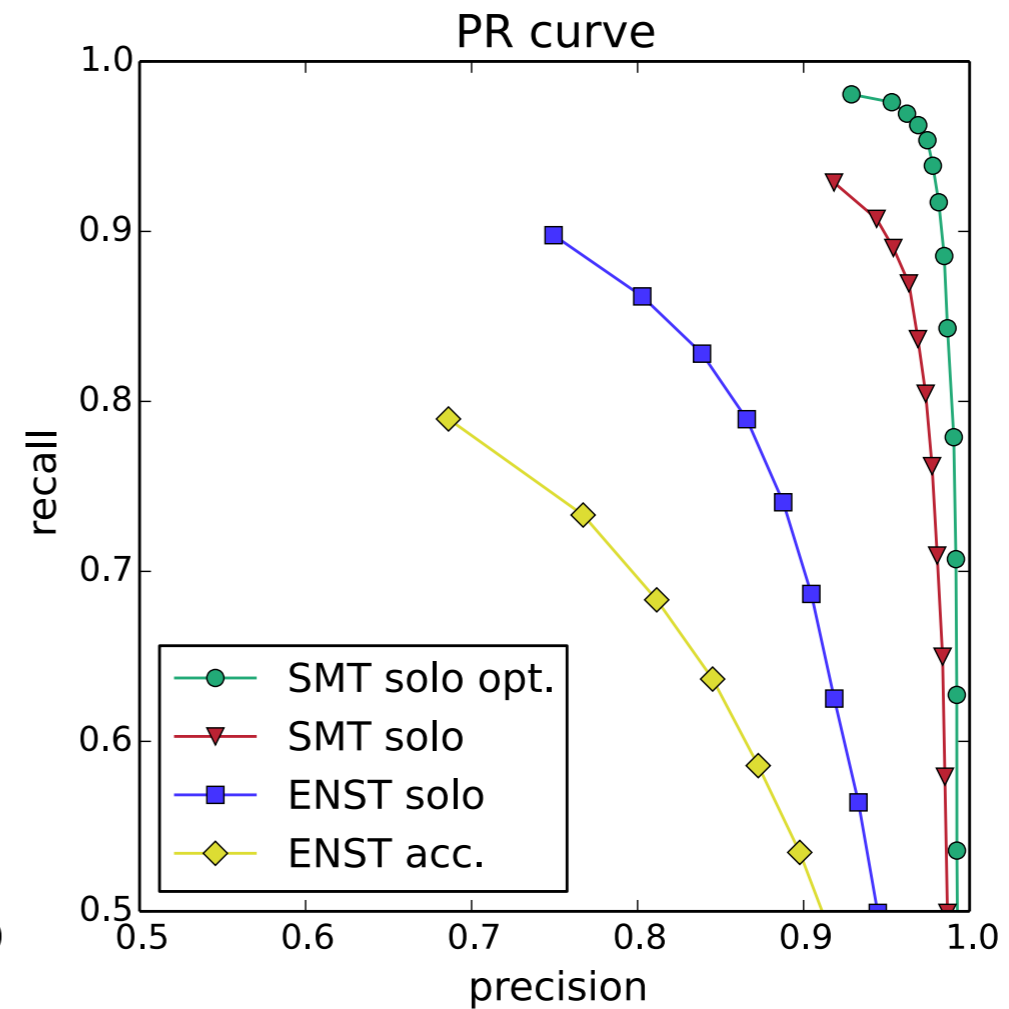
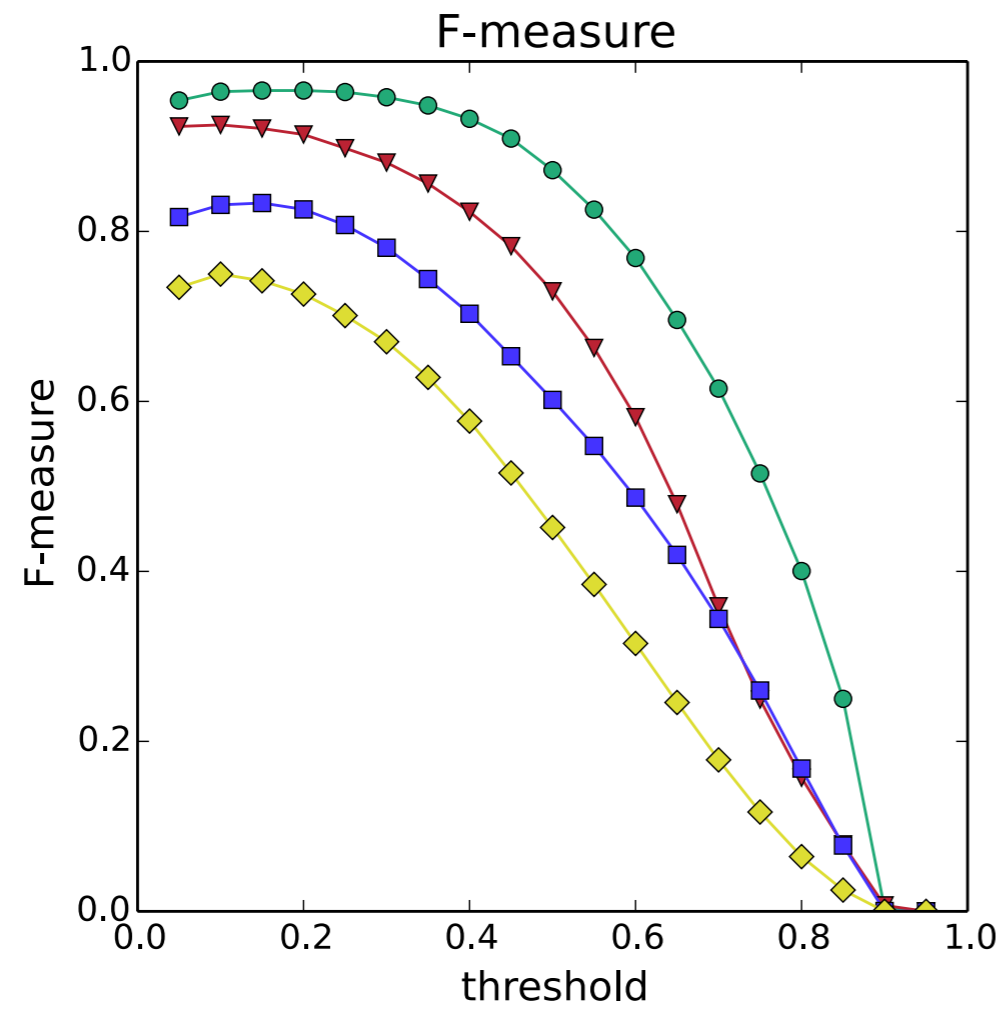
# RESULTS

Method	SMT opt.	SMT solo	ENST solo	ENST acc.
<b>NMF-SAB</b> [Dittmar and Gärtner 2014]	95.0	—	—	—
<b>PFNMF</b> [Wu and Lerch 2015]	—	81.6	77.9	72.2
<b>HMM</b> [Paulus and Klapuri 2009]	—	—	81.5	74.7
<b>BDRNN</b> [Southall et al. 2016]	96.1	83.3	73.2	66.9
<b><i>tsRNN</i></b>	<b>96.6</b>	<b>92.5</b>	<b>83.3</b>	<b>75.0</b>
	$\delta = 0.10$	$\delta = 0.15$	$\delta = 0.15$	$\delta = 0.10$

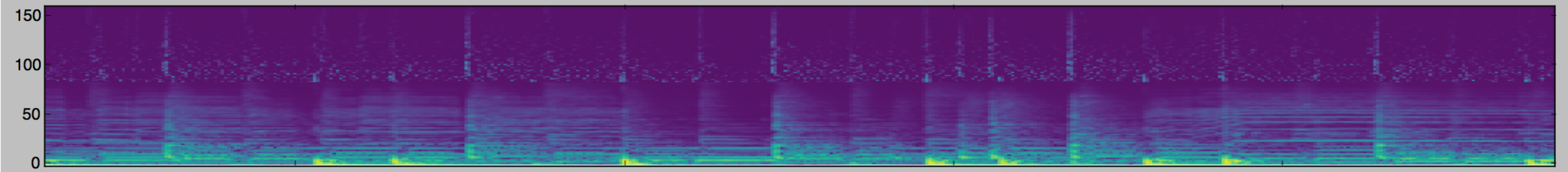
# RESULTS

Method	SMT opt.	SMT solo	ENST solo	ENST acc.
<b>NMF-SAB</b> [Dittmar and Gärtner 2014]	95.0	—	—	—
<b>PFNMF</b> [Wu and Lerch 2015]	—	81.6	77.9	72.2
<b>HMM</b> [Paulus and Klapuri 2009]	—	—	81.5	74.7
<b>BDRNN</b> [Southall et al. 2016]	96.1	83.3	73.2	66.9
<b><i>tsRNN</i></b>	<b>96.6</b>	<b>92.5</b>	<b>83.3</b>	<b>75.0</b>
	$\delta = 0.10$	$\delta = 0.15$	$\delta = 0.15$	$\delta = 0.10$

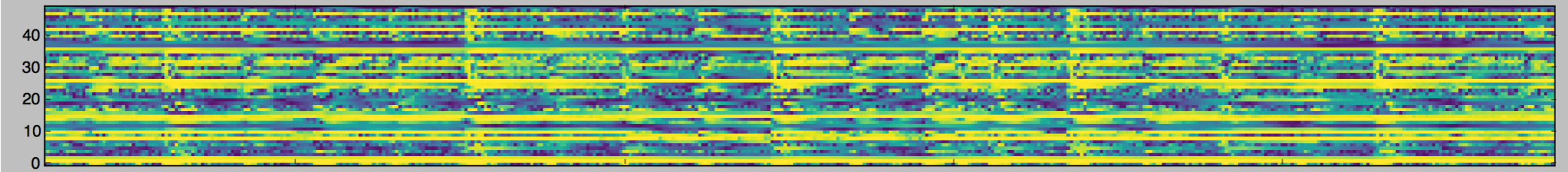
# RESULTS



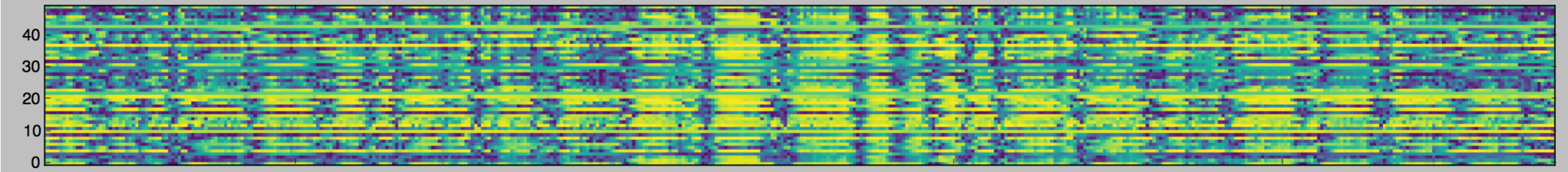
Input



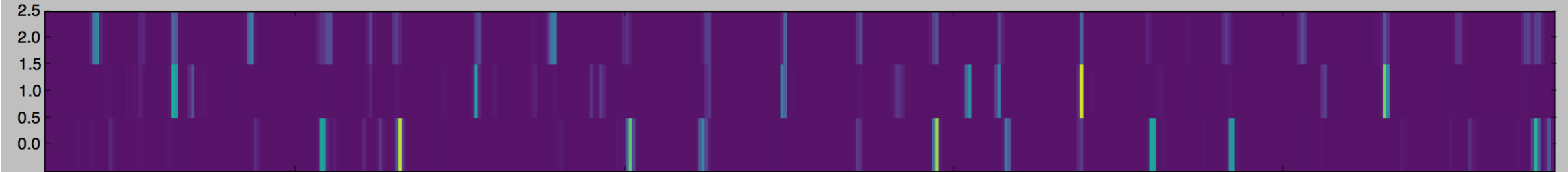
GRU1



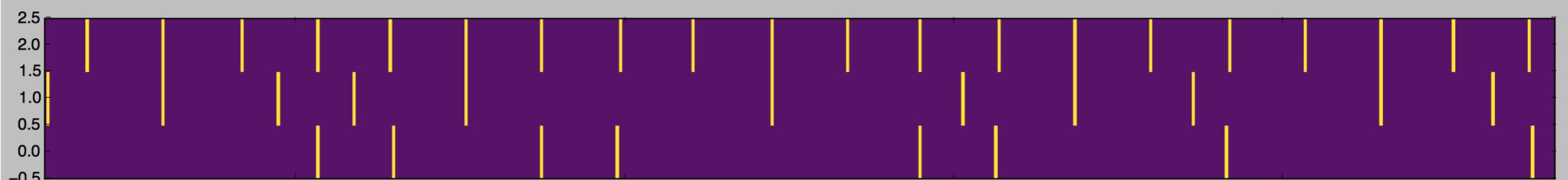
GRU2



Output



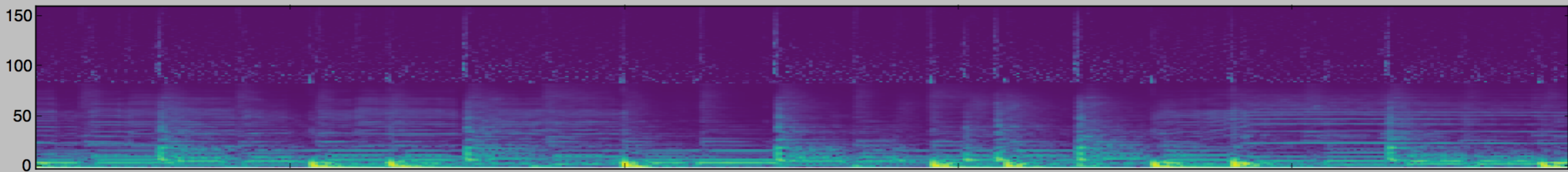
Targets



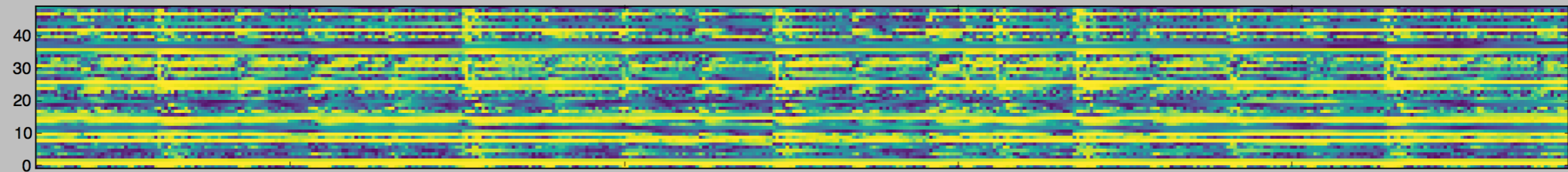
2700 2800 2900 3000

Time ->

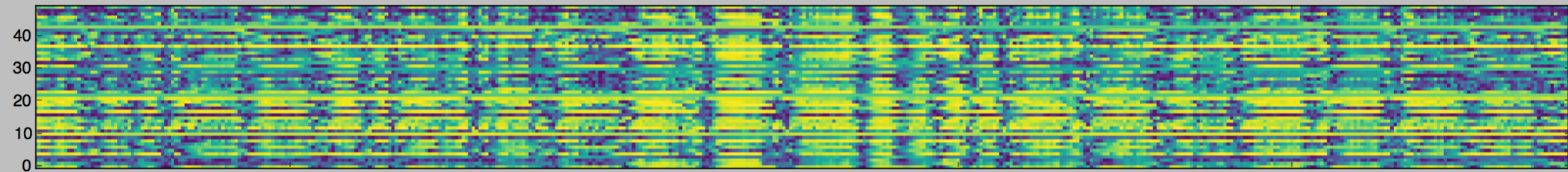
Input



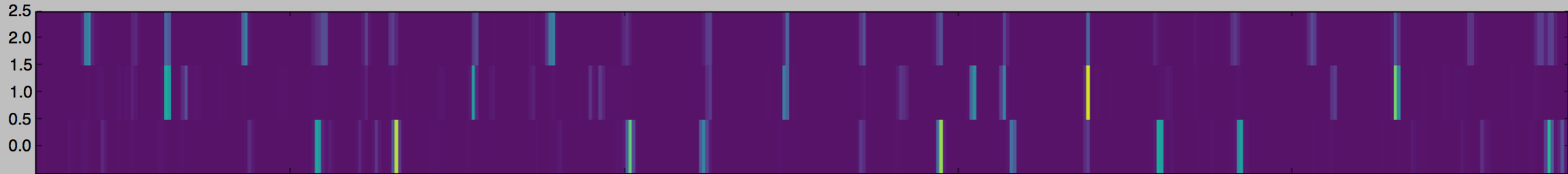
GRU1



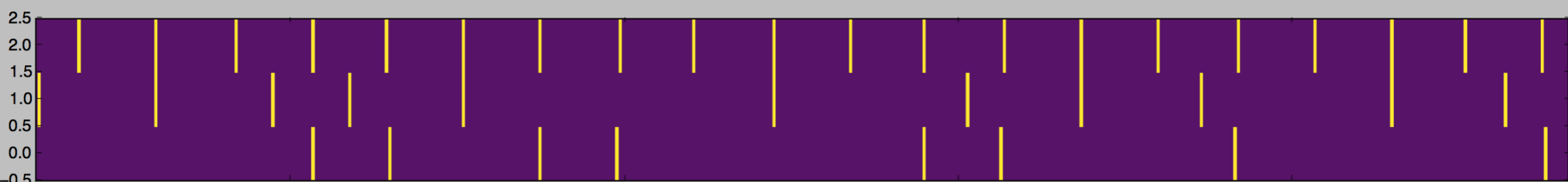
GRU2



Output



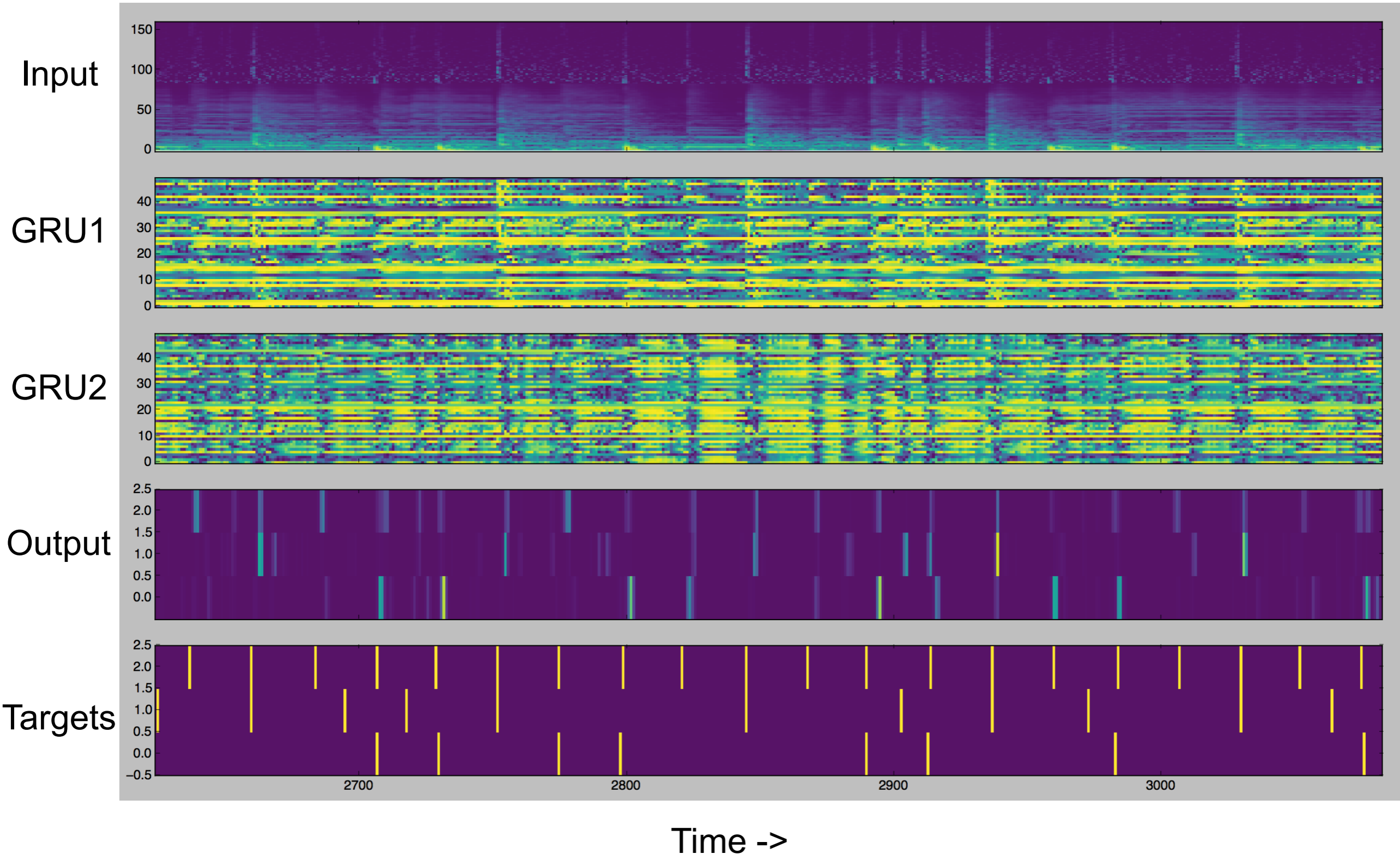
Targets



2700 2800 2900 3000

Time ->







# CONCLUSIONS

- Towards a **generic end-to-end acoustic model for drum detection** using **RNNs**
- **Data augmentation** greatly improves generalization
- **Weighting loss functions** helps to improve detection of difficult instruments
- RNNs with label **time shift** perform equal to BDRNN
- Simple RNN architecture **performs better or similarly well** as handcrafted techniques  
while using a **smaller tolerance** window (20ms)