# INCREASING DRUM TRANSCRIPTION VOCABULARY USING DATA SYNTHESIS

Mark Cartwright and Juan Pablo Bello

Music and Audio Research Laboratory

New York University
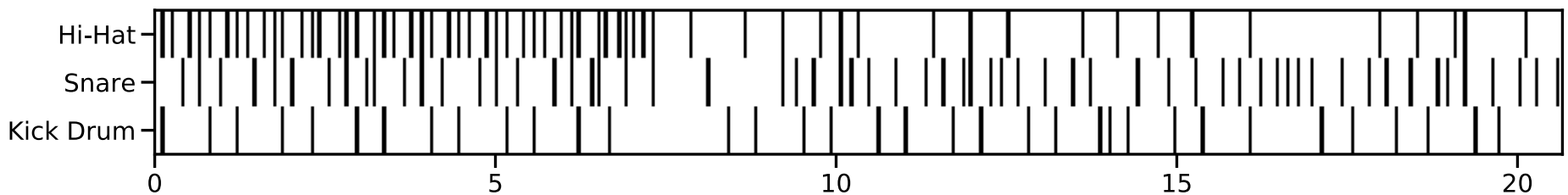
http://www.markcartwright.com

1

# Problem

Most Automatic Drum Transcription (ADT) algorithms are limited to simply onset times of 3 classes:

1. Bass Drum (BD)
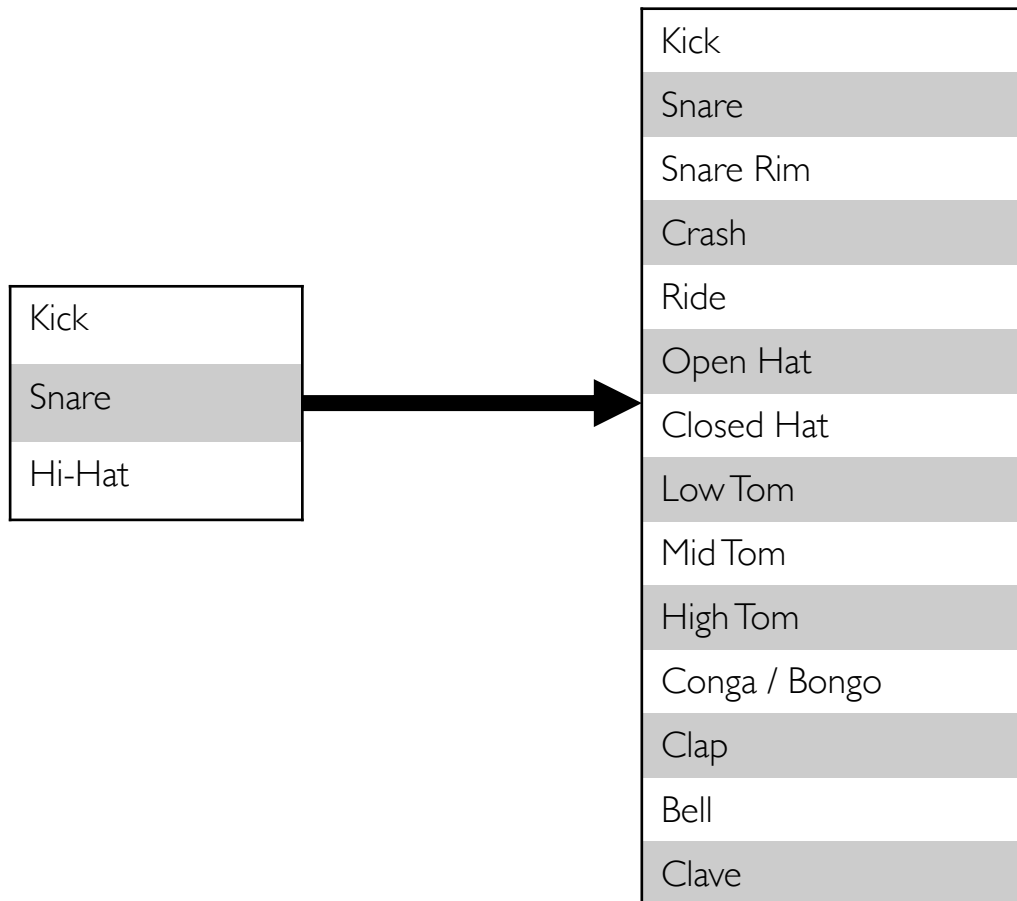2. Snare Drum (SD)
3. Hi-Hat (HH)

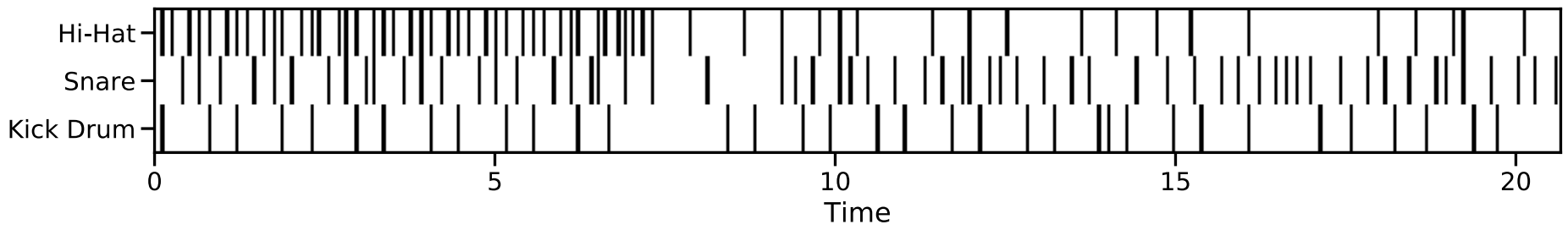Example:

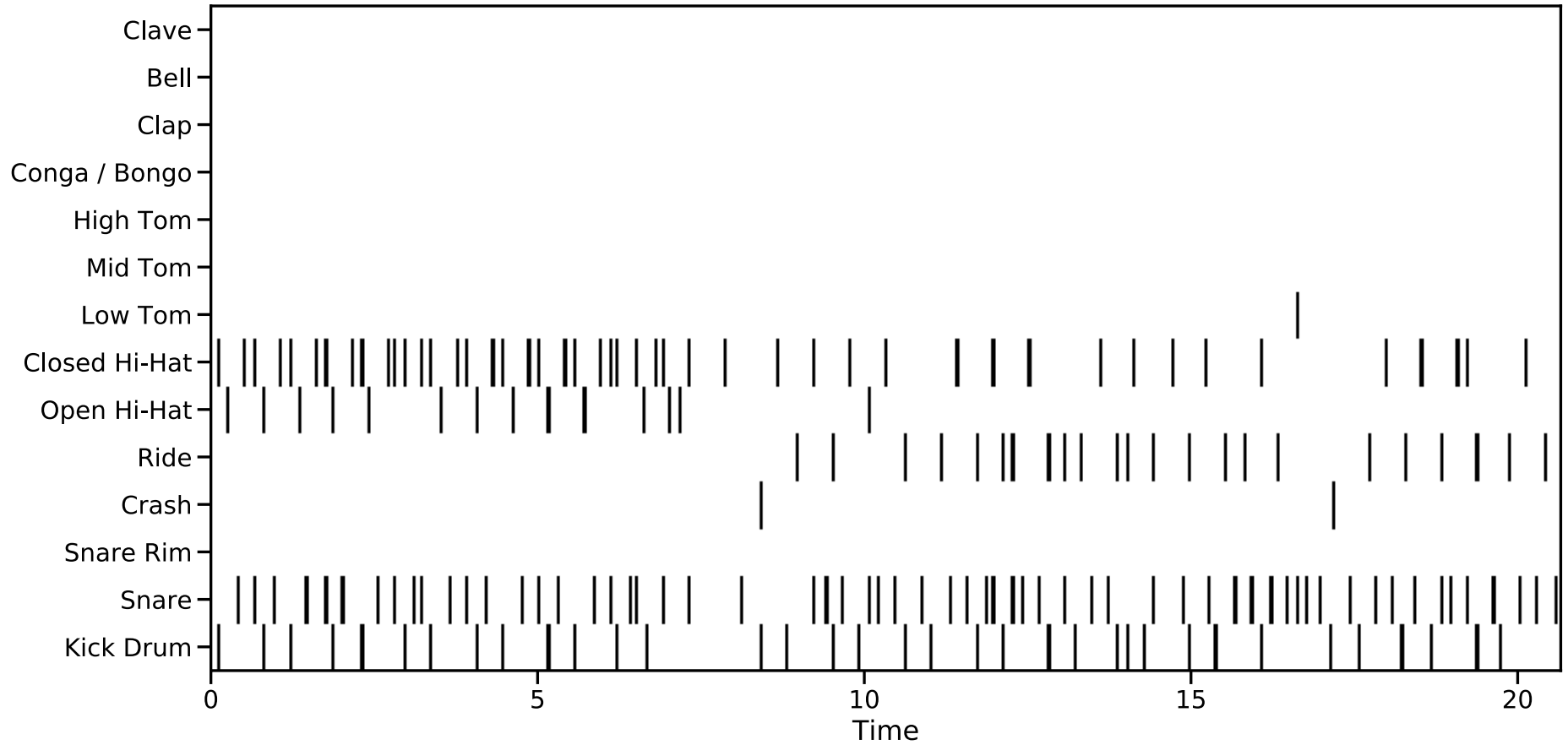Original 🔊))          3-class re-synthesis: 🔊))

# Goal

- Increase vocabulary to onset times of 14 drum classes:



| Kick |
| Snare |
| Hi-Hat |

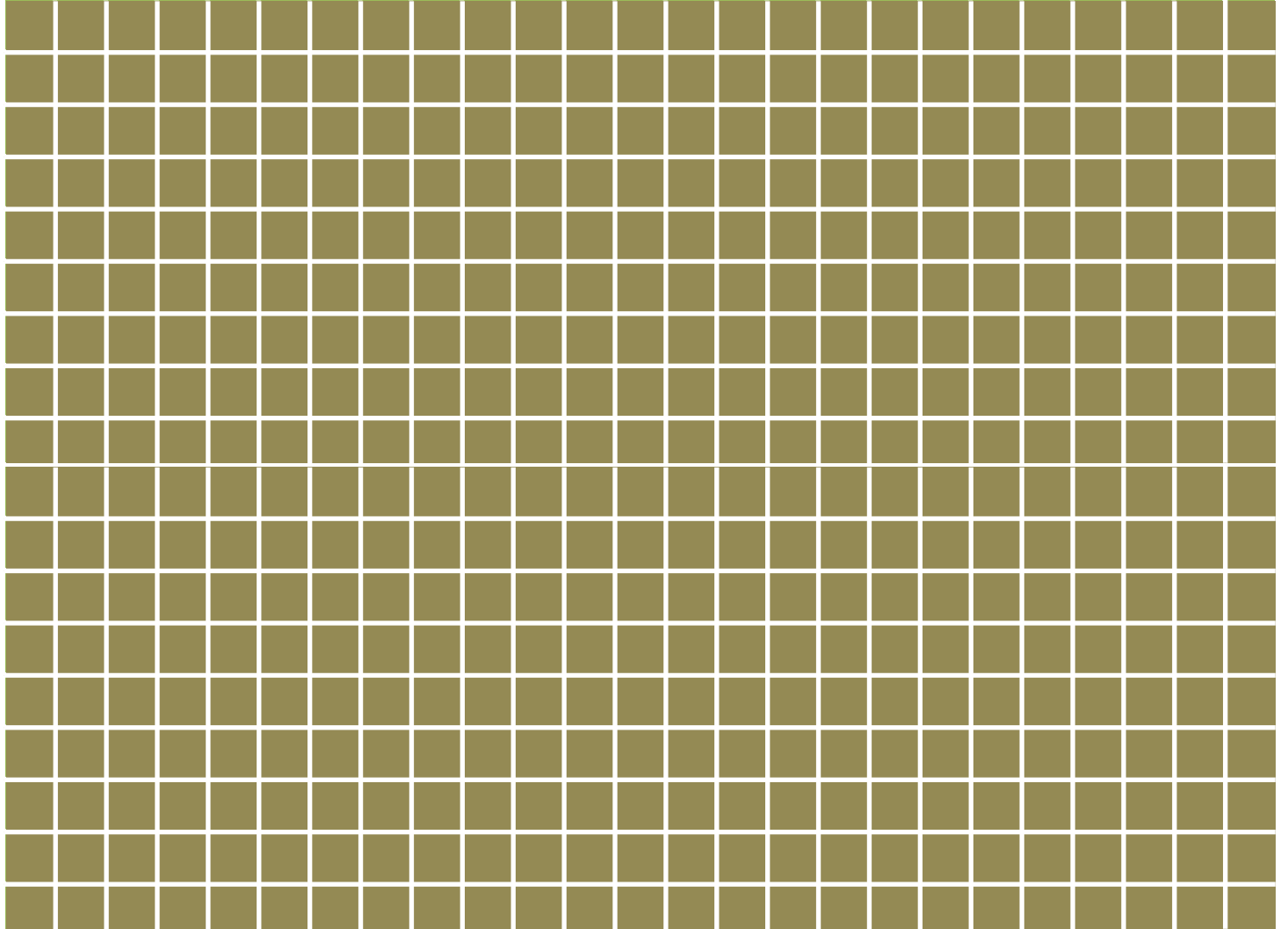| Kick |
| Snare |
| Snare Rim |
| Crash |
| Ride |
| Open Hat |
| Closed Hat |
| Low Tom |
| Mid Tom |
| High Tom |
| Conga / Bongo |
| Clap |
| Bell |
| Clave |

# Goal

Goal

# ADT Datasets

~33,000 onsets / 2 hours of 3-class data
(RBMA, IDMT/SMT)

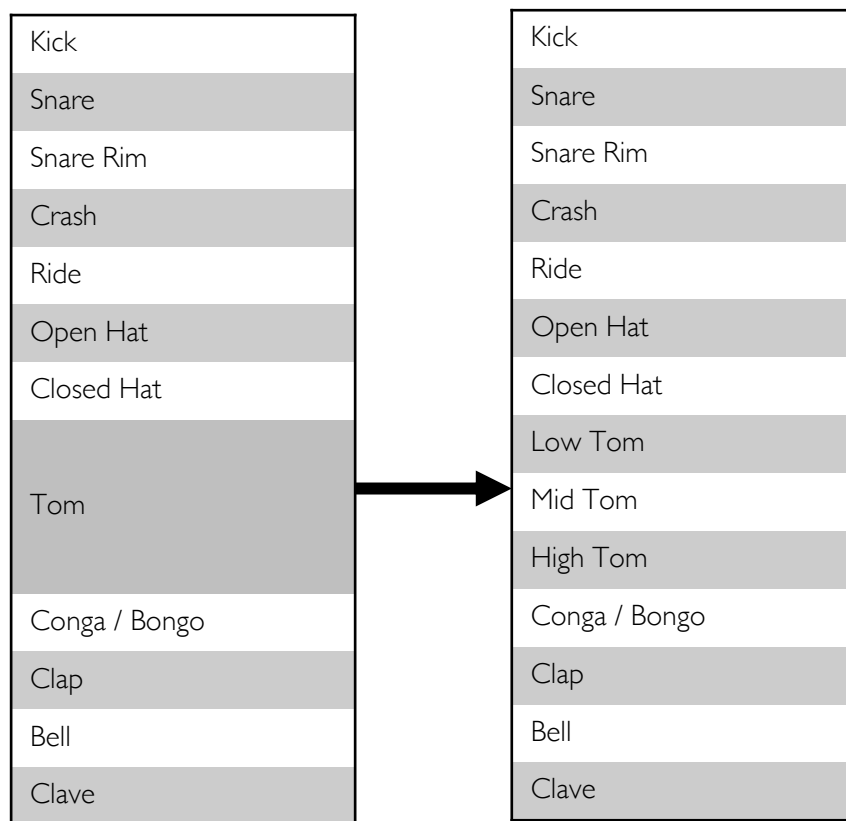~33,000 onsets / ~1.5 hours of > 3-class data
(ENST/MDB)

# ADT Datasets

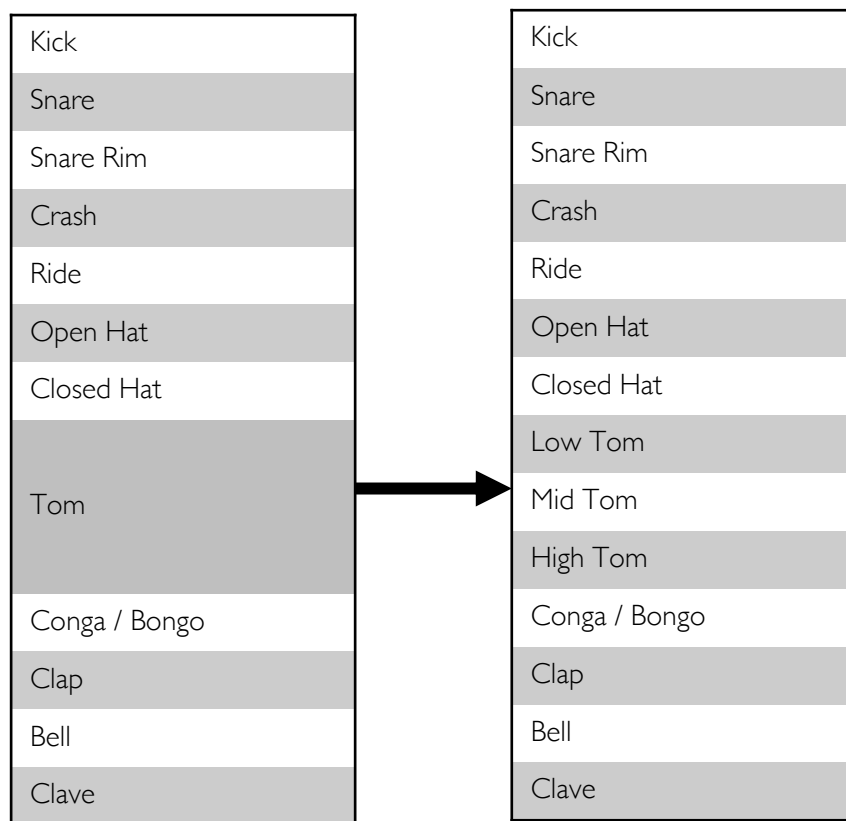# Building the S********* Drum Dataset (SDDS)
## *The Sounds*

- Labeled samples from 8 drum sample libraries
- Split toms into low, mid, high based on pitch, spectral centroid
- Split dataset into 4000 train / 2000 test samples

| Kick |
|---|
| Snare |
| Snare Rim |
| Crash |
| Ride |
| Open Hat |
| Closed Hat |
| Tom |
| Conga / Bongo |
| Clap |
| Bell |
| Clave |

→

| Kick |
|---|
| Snare |
| Snare Rim |
| Crash |
| Ride |
| Open Hat |
| Closed Hat |
| Low Tom |
| Mid Tom |
| High Tom |
| Conga / Bongo |
| Clap |
| Bell |
| Clave |

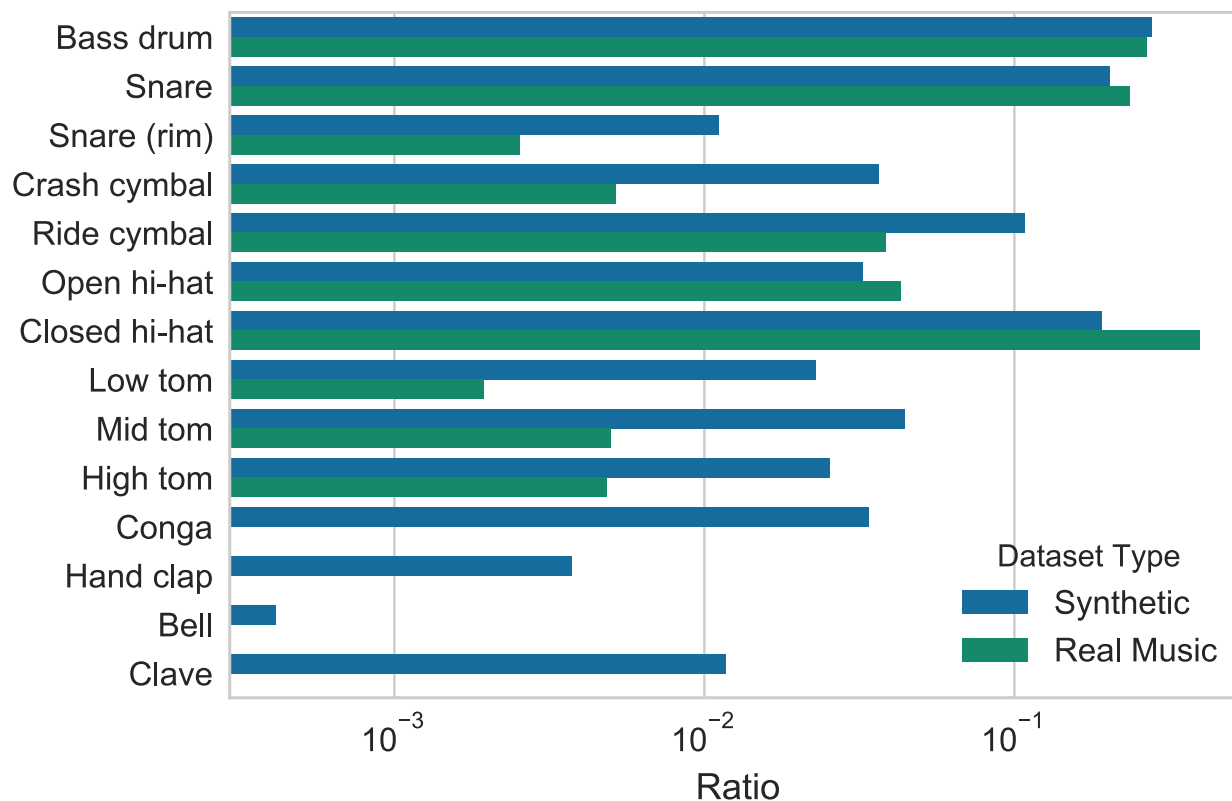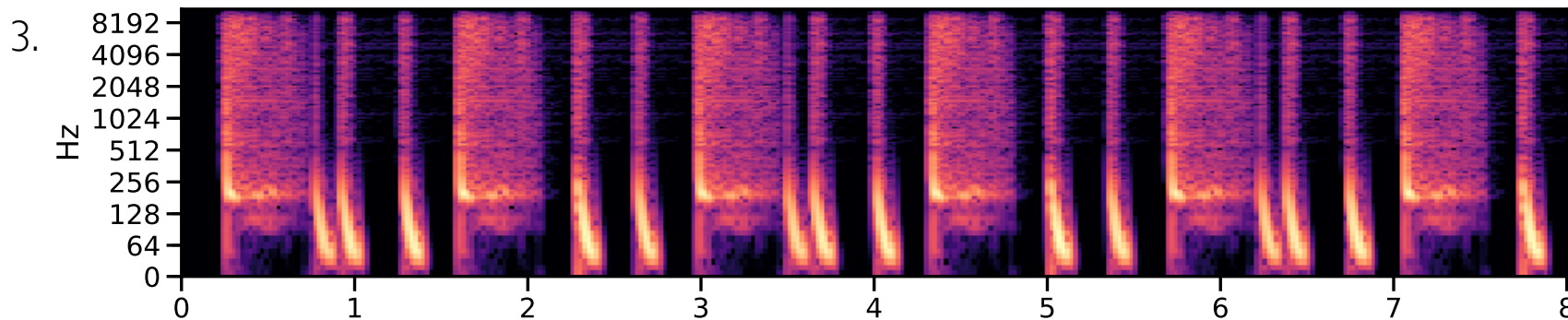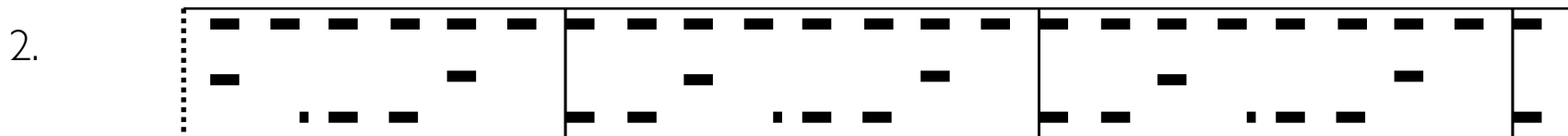# Building the Synthetic Drum Dataset (SDDS)
## *The Sounds*

- Labeled samples from 8 drum sample libraries
- Split toms into low, mid, high based on pitch, spectral centroid
- Split dataset into 4000 train / 2000 test samples

# Building the Synthetic Drum Dataset (SDDS)
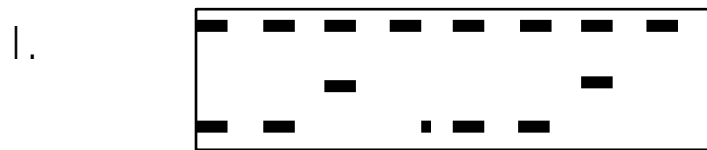## *The Sounds*

- Labeled samples from 8 drum sample libraries
- Split toms into low, mid, high based on pitch, spectral centroid
- Split dataset into 4000 train / 2000 test samples

# Building the Synthetic Drum Dataset (SDDS)
## *The Rhythms*

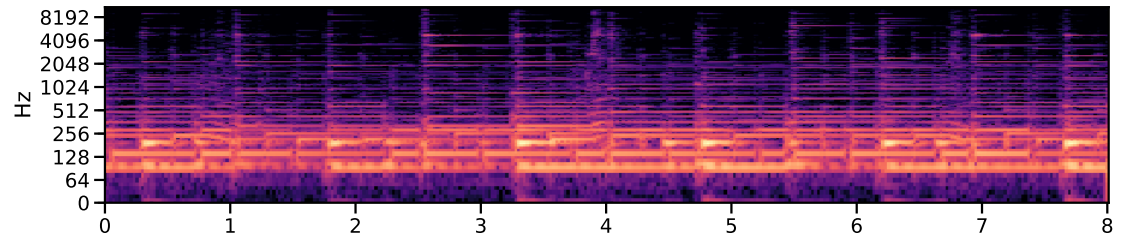- 60k measures of percussion MIDI files (50k train / 5k test / 5k validate)

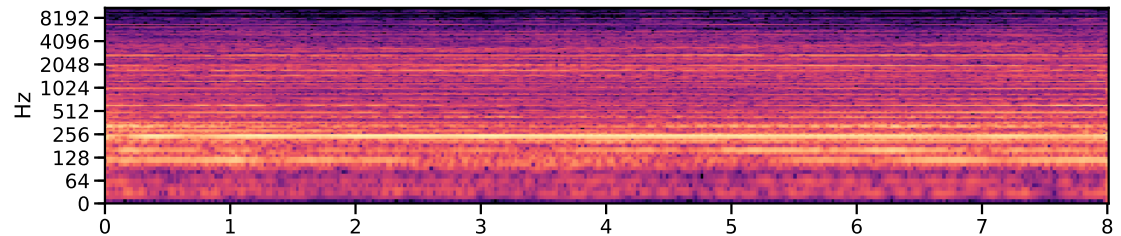# Building the Synthetic Drum Dataset (SDDS)
## *Augmentation and "Accompaniment"*

- Augment to 210k (200k train / 5k test / 5k validate)  w/ small pitch shifts, added pink noise, and "*harmonic noise* accompaniment":
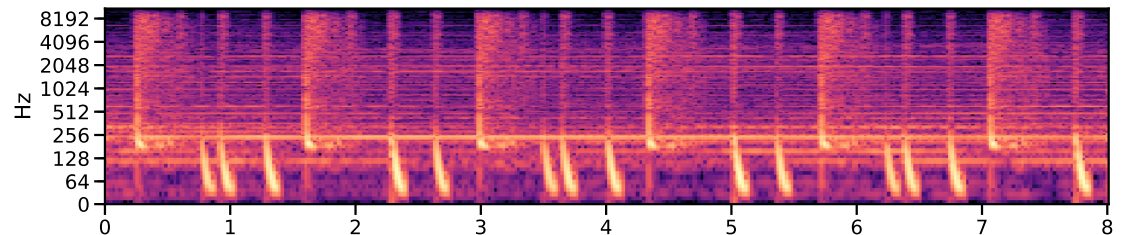
- Solo harmonic instrument recording



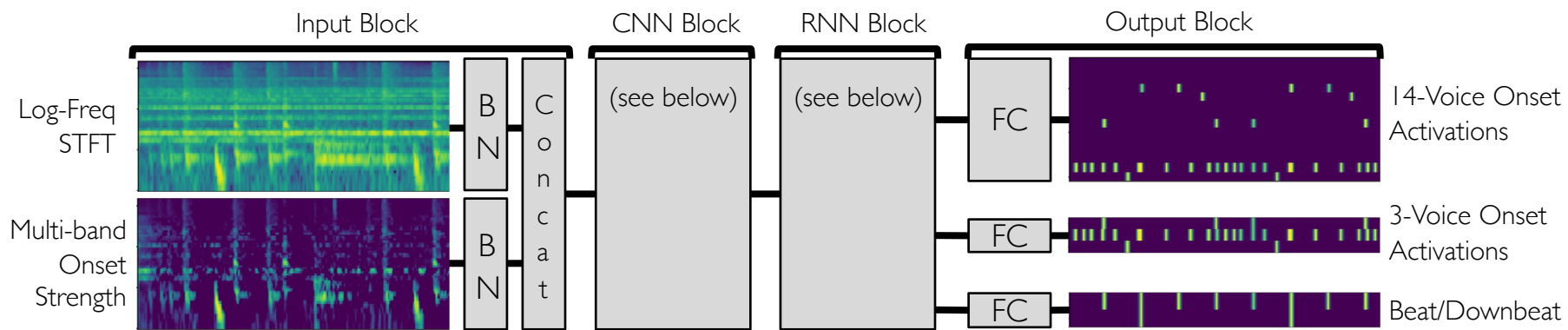- "Smear" in time (fwd / bkwd reverb)



- Mix with drum track

# Combined Datasets

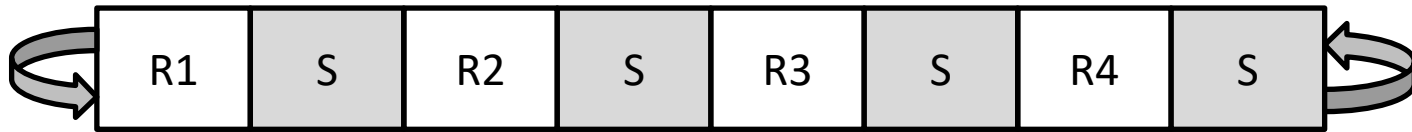| | RBMA | SMT | ENST | MDB | SDDS |
|---|---|---|---|---|---|
| Hours | 1.67 | 0.51 | 1.28 | 0.23 | 467 |
| Accomp. | X | X | | X | sort of |
| 14-voice Onsets | | | X | X | X |
| 3-voice Onsets | X | X | X | X | X |
| Beats | X | | | | X |

# Model



| Input Features |
| --- |
| 22050 sampling rate |
| **Log-magnitude, log-frequency STFT:**<br>1024 samples<br>64 bands (8 octaves, 8 bands per octave)<br>40 Hz – 10kHz<br>0.01 sec hop |
| **Multi-band onset strength:**<br>Computed from log-f STFT<br>Half-rectified difference between current<br>frame and mean of past 200 ms |

| CNN Block |
| --- |
| 32 (3x3) Conv |
| 32 (3x3) Conv |
| Batch Norm |
| ReLU |
| 30% Dropout |
| 64 (3x3) Conv |
| 64 (3x3) Conv |
| BatchNorm |
| ReLu |
| 30% Dropout |
| 64 (1x64) Conv |
| BatchNorm |
| ReLU |

| RNN Block |
| --- |
| (-6:+6) Context Windowing |
| 64-unit BLSTM |
| 64-unit BLSTM |
| 64-unit BLSTM |

# Training with Heterogeneous Outputs

- Mask outputs not in use for each example
- Use round-robin sampling with Pescador[1] to ensure all outputs used in each mini-batch of 8:



- Minimize weighted combination of binary cross-entropy losses based using weights computed by activation entropy
- 3-fold CV splits for the small real music small datasets 25% validation / 75% testing in each split
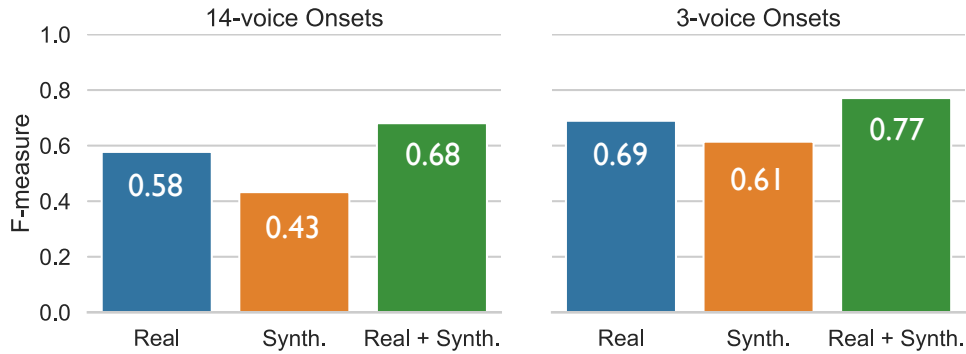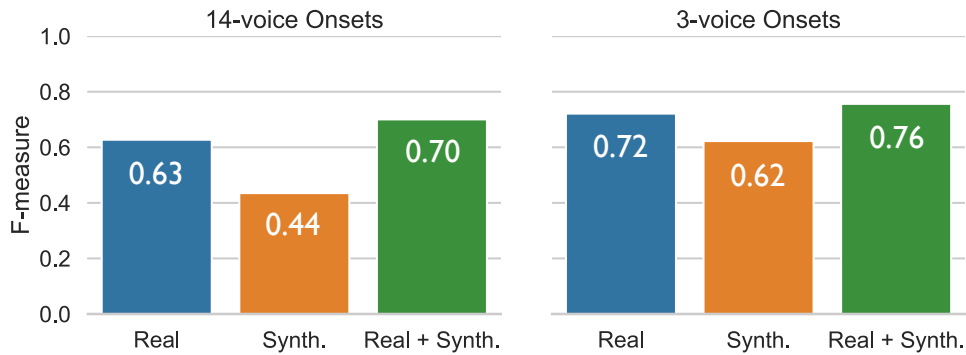
1. https://github.com/pescadores/pescador

# Experiments

Variables:
- Training data
  - Real music (RBMA, SMT, ENST, MDB)
  - Synthetic (SDDS)
  - Recorded + Synthetic

- Model capacity
  - "Small" (as described)
  - Large (more conv filters, more BLSTM units)

- Outputs
  - Multi-task
  - Single-task (w/ limited data)

- Class-weighting
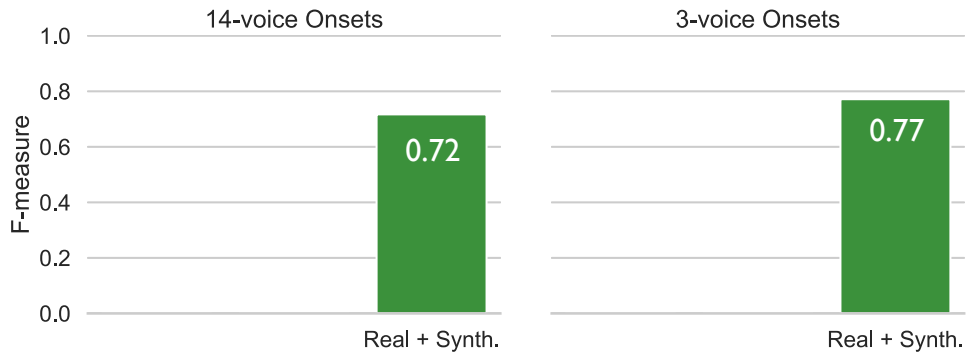  - No weighting
  - Weighted by activation entropy

# Results

# Results

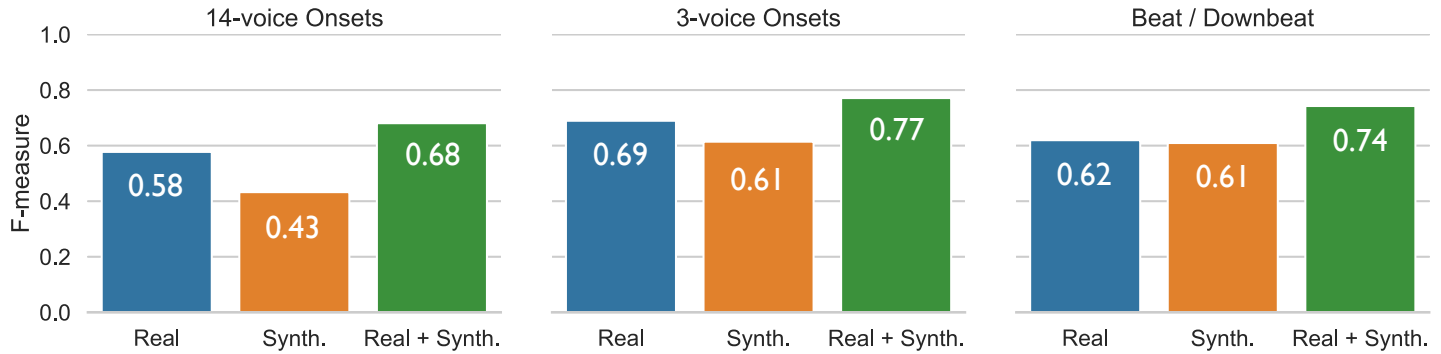Training: Real | Evaluation: Real

# Results

Training: Real | Evaluation: Real

Macro F1: 0.20

Training: Real + Synth | Evaluation: Real

Macro F1: 0.35

Weighted - Training: Real + Synth | Evaluation: Real

Macro F1: 0.61
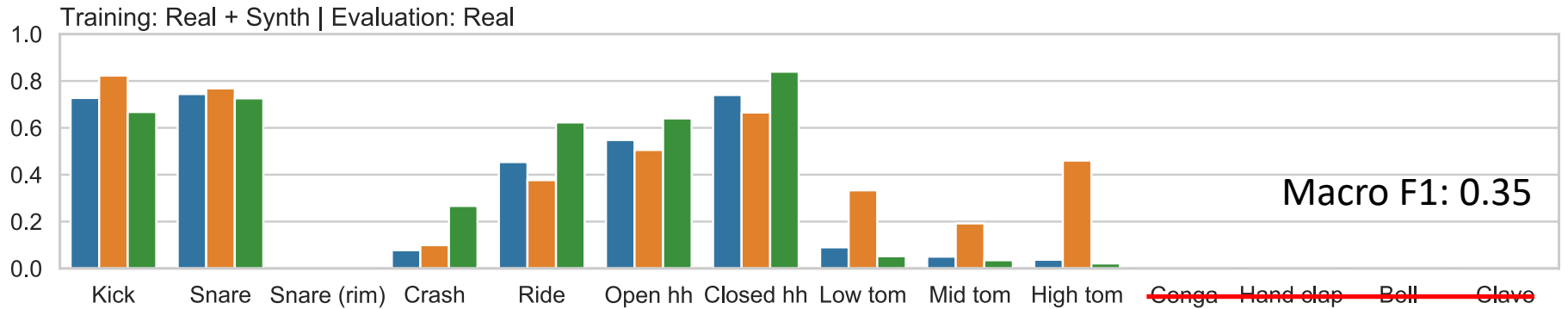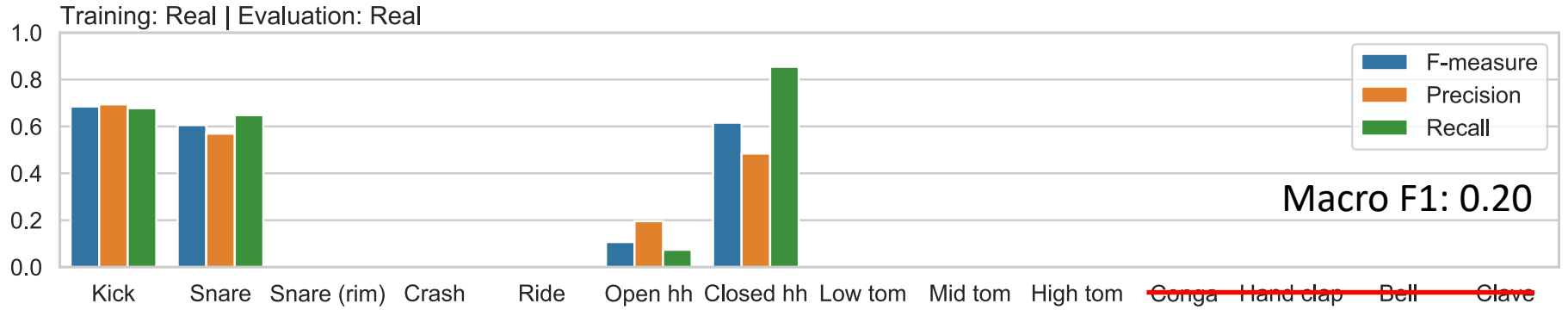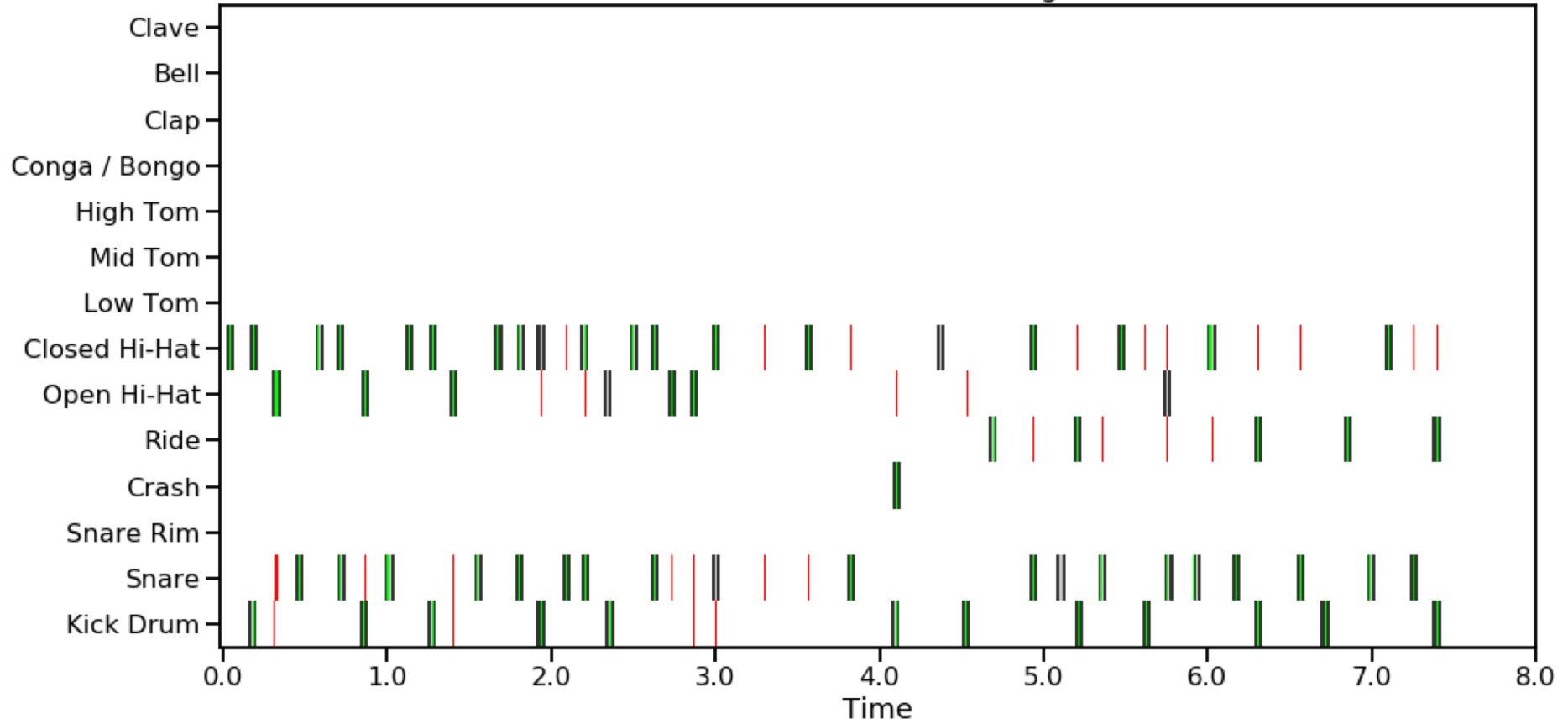
# Results



Small / Multi-task / Class-weighted

# Conclusion

- Lots of mediocre synthetic data can improve performance on both overall performance and uncommon classes

- However, it must be used in conjunction with real music data

- Multi-task learning doesn't seem to help for large-vocab transcription, but does help in the auxiliary task of downbeat/beat tracking

Download the trained models at
https://github.com/mcartwright/dafx2018_adt