

# DrummerNet – Deep Unsupervised Drum Transcription



Keunwoo Choi and Kyunghyun Cho  
Spotify and New York University  
keunwooc@spotify.com

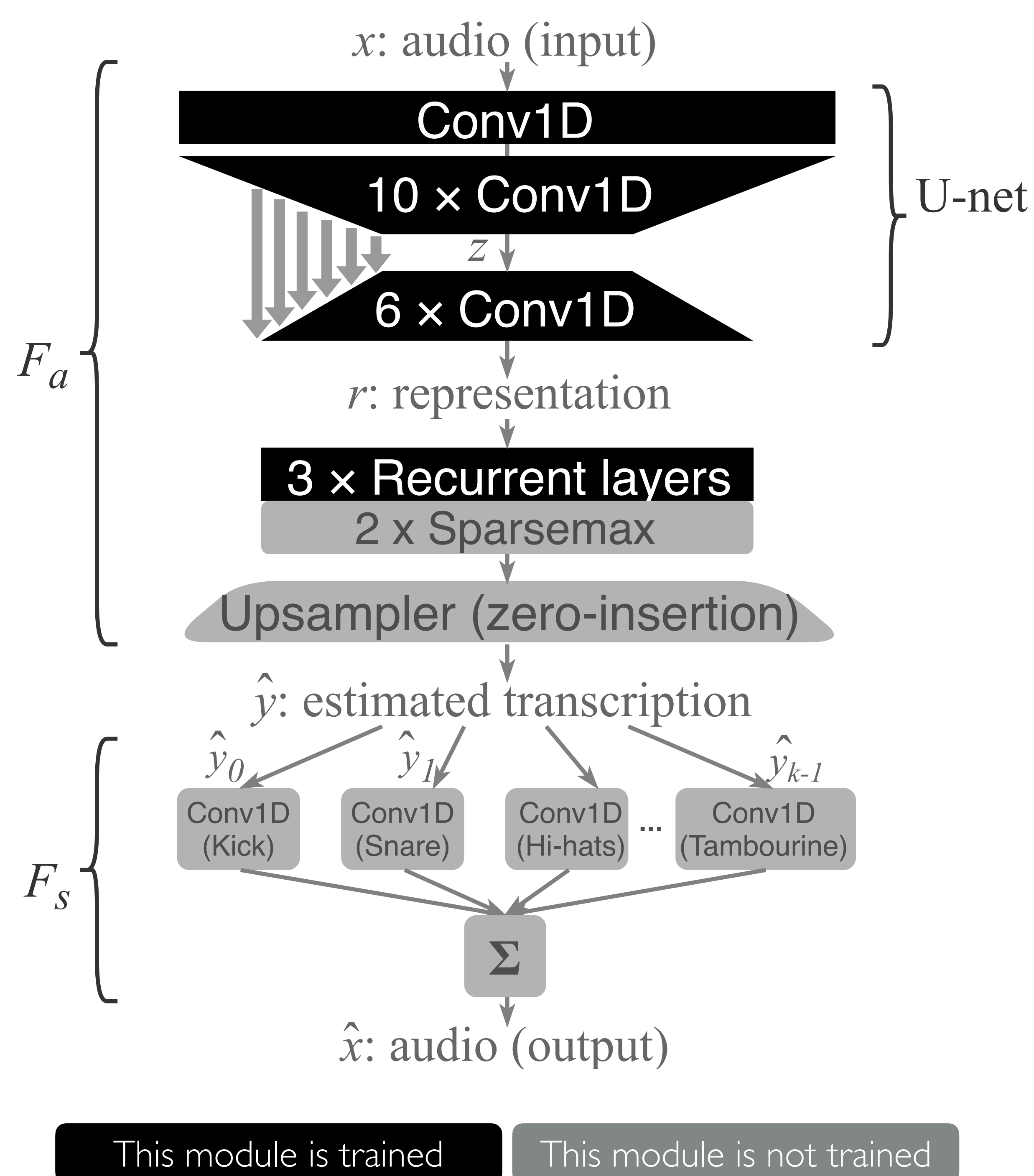
Code + Paper



## 1. Introduction

- ▶ def(Drum Transcription): drum stem  $\rightarrow$  music score with kick/snare/hihat
- ▶ Unsupervised transcription: transcribe using audio only!
  - ▶ Because it's expensive to get an large annotated dataset
    - ▶ Hence supervised learning approaches don't generalize well
      - ▶ i.e., fails with unseen kind of drum sound
- ▶ Any prior work?
  - ▶ Few non-deep unsupervised (drum, or others) transcription systems

## 3. DrummerNet



- ▶ **Analysis system,  $F_a$** : Transcriber
  - ▶ **Unet** learns a representation  $r$  based on the input audio drum stem  $x$
  - ▶ **Recurrent layers** learn to transcribe along time/instruments
  - ▶ **Sparsemax layers** output \*really\* sparse activation (lots of zeros)
  - ▶ **Upsampler** outputs transcription  $y_{hat}$ 
    - ▶ Zeros are inserted to compensate downsampling effect of U-net
- ▶ **Synthesis system,  $F_s$** : Synthesizer
  - ▶ **Conv1D layers** synthesize each drum track by convolving  $\{\text{impulses } (y_{hat\_k})\} * \{\text{drum component waveform\_k}\}$
  - ▶  **$\Sigma$**  sums up each track  $x_{hat\_k}$
- ▶ **Training**
  - ▶ Minimize reconstruction error ( $x$  vs  $x_{hat}$ )

## 2. How is it possible?

- ▶ Imagine a structure like this:



Q. If we train Blackbox so that  $x_{hat}$  becomes similar to  $x$ , what would Blackbox do?

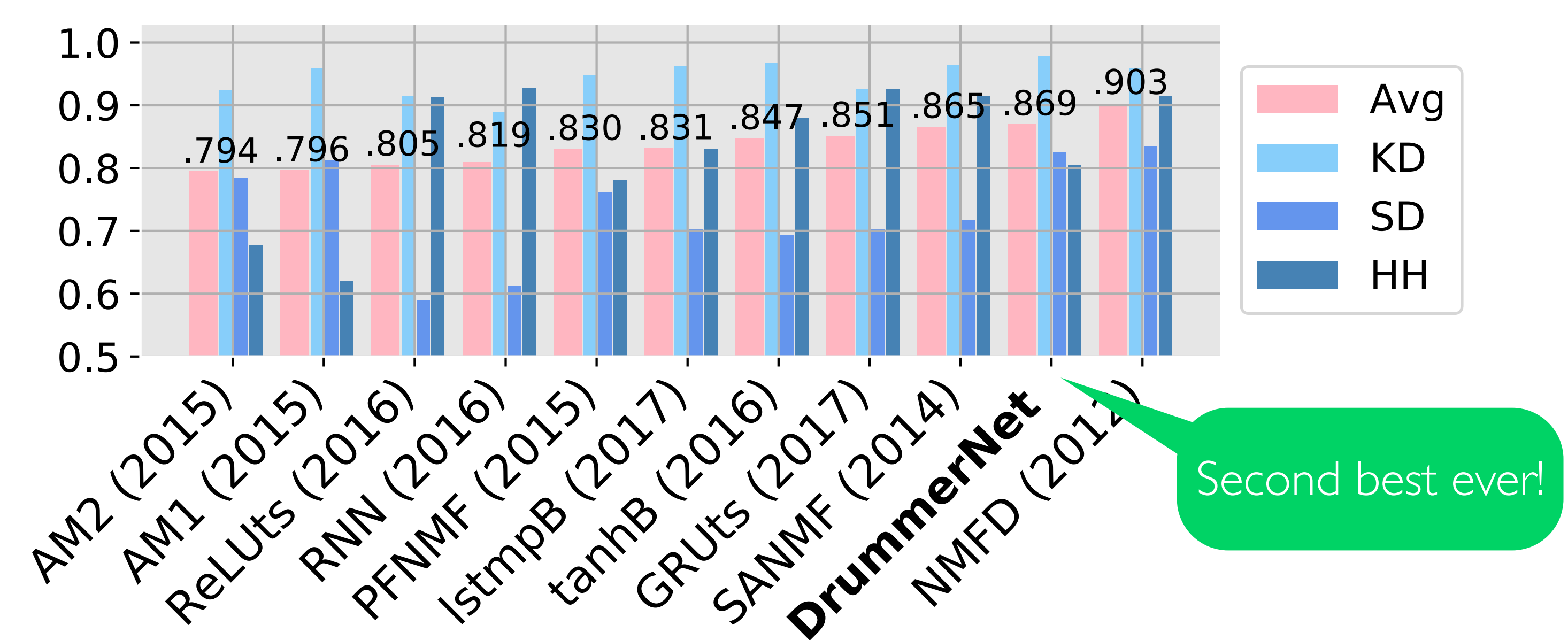
*\*\*drum rolls\*\**

A. Drum Transcription!

## 4. Results

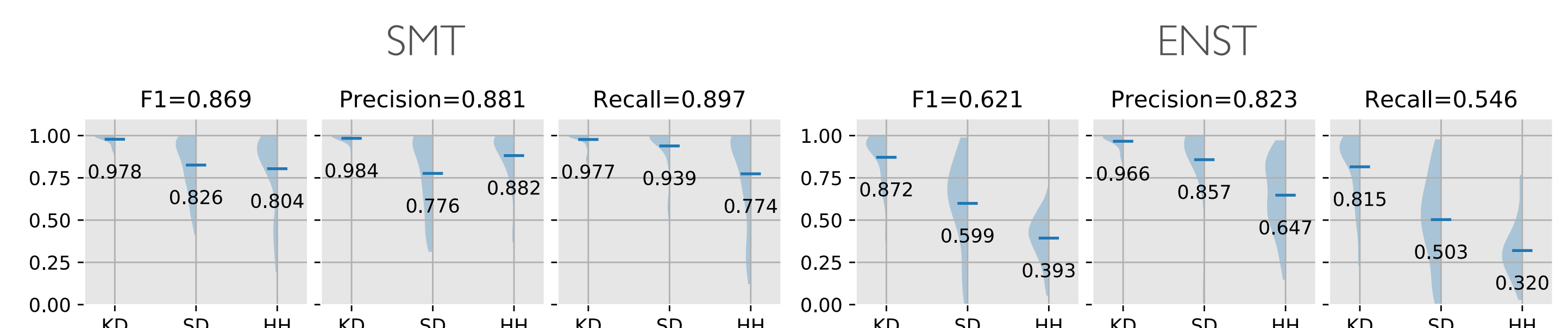
- ▶ F1 scores on a dataset ('SMT') compared to others

- ▶ DrummerNet uses much larger training dataset and it leads to better generalizability

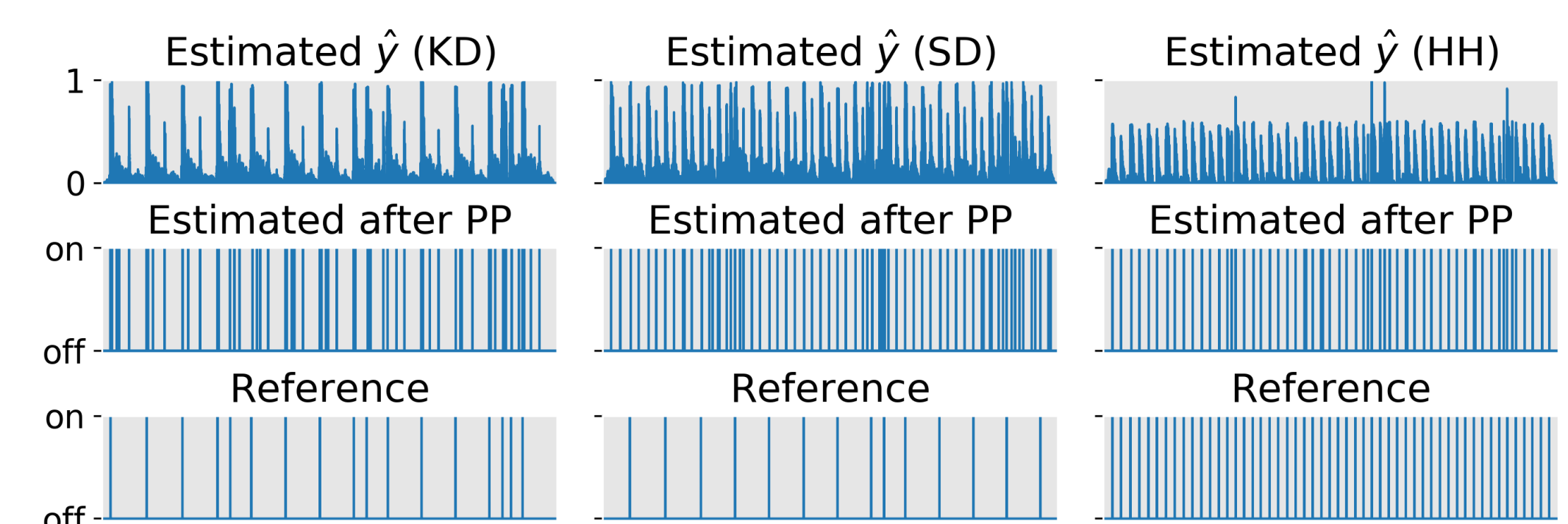


- ▶ SMT vs ENST dataset

- ▶ Because the reconstruction error is measured by audio domain, there's linearity between {transcription, audio} + We mix-use the concepts of velocity vs probability  $\rightarrow$  we get false positives if
  - i) different drum components sound (sort of) similar
  - ii) the play is nuanced (**var(velocities)** is large)  $\rightarrow$  e.g., ENST dataset



- ▶ DrummerNet fails like this when it fails



- ▶ More results and discussion including ablation study are in the paper!

## 5. Future work/Code/Paper

Future work

Universal unsupervised transcription!  
RL + Unsupervised transcription

code/paper

<https://github.com/keunwoochoi/DrummerNet>

