

Automatic Polyphonic Music Transcription

ECE 208/408 *Spring*
2025

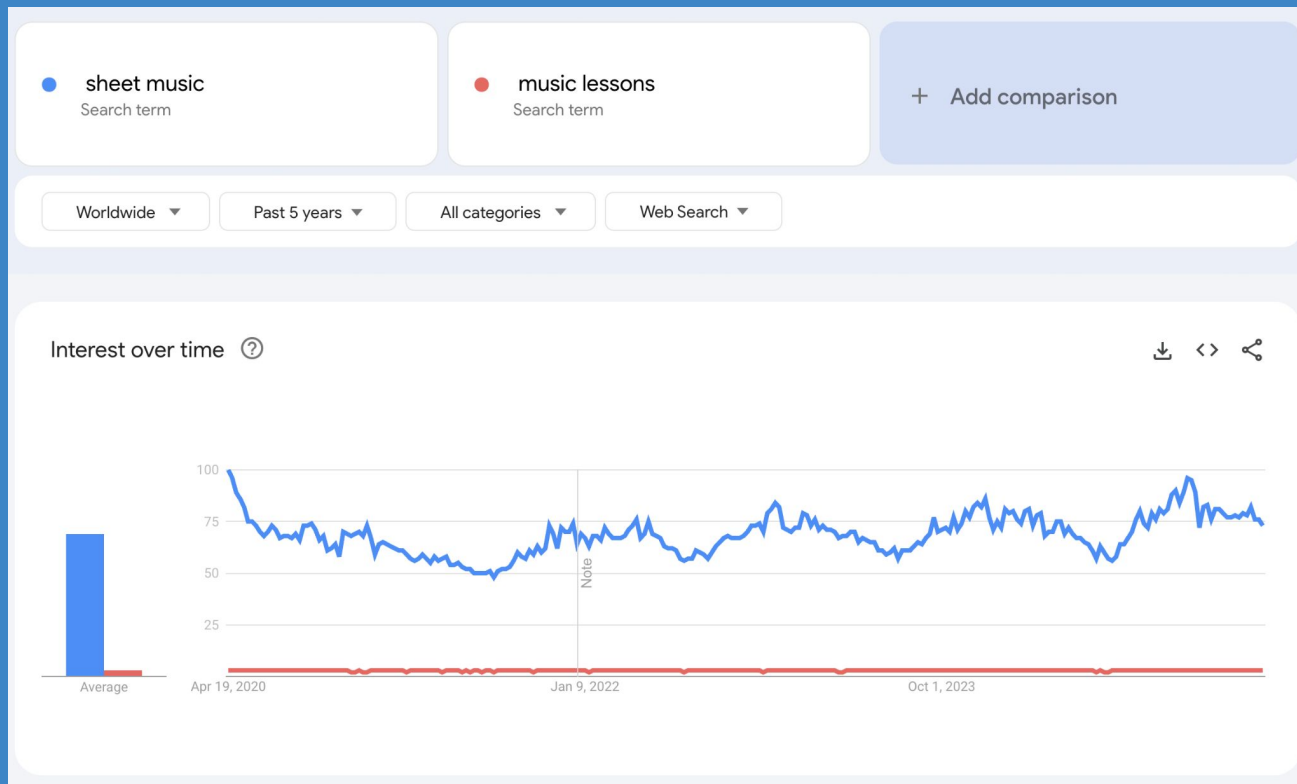
R. Faran, Y. Wu, S. Xu

01

Motivation

AUDIO TRANSCRIPTION

Continued need for accessible music



MUSIC TRANSCRIPTION

Lack of dedicated transcribers

- Very few professional musicians can make a living off of transcription alone

Years of training required

- Often overlap with composers or arrangers

Music Publishers are constantly looking for better and cheaper methods

Current models inaccurate, especially for multiple notes



"Although some may manage to make a living working exclusively as transcribers, for most, this job is just one of many ways for highly trained musicians to use their skill sets in an entrepreneurial fashion"

<https://www.berklee.edu/careers/roles/transcriber>

Why Classical Music?

High cost of entry/lessons limits accessibility

- \$60/hr or even higher!

Very old genre

- Centuries worth of pieces and decades of recordings

Common curriculum of study for musicians and schools

- Education and scholarship

Accessibility of dataset: most pop music is copyrighted



02

Setup

Music Transcription Steps

01 DATASET CURATION - Faran

- .wav files included
- Manually Transcribed and verified by trained musicians
- Variety of instrumentation
- Split training and test data

02 PRE-PROCESSING - Faran

- Load dataset
- STFT and Chroma data
- Load training labels (.csv)
- Infer Timing for labels

03 TRAIN A MODEL - Wu, Xu

- CNN
- HMM
- Predictive Modeling

04 TUNE HYPERPARAMETERS- Wu

- Number of layers
- Batch and hop sizes
- Layer Connectivity

05 OUTPUT EVALUATION - Faran, Xu

- .wav Input
- MIDI Output
- Accuracy

03

Our Model

PRE-PROCESSING

Loading the Dataset and Labels

-reorganizing the csv, extracting labels and link it to the corresponding wav file

Only uses start_time, end_time and notes

-load .wav as mono 1D Numpy array
-Downsample the audio to 22050 Hz

start_time	end_time	instrument	note	start_beat	end_beat	note_value
9182	90078	43	53	4.0	1.5	Dotted Quarter
9182	33758	42	65	4.0	0.5	Eighth
9182	62430	1	69	4.0	1.0	Quarter
9182	202206	44	41	4.0	3.5	Whole
9182	62430	1	81	4.0	1.0	Quarter
33758	62430	42	60	4.5	0.5	Eighth
62430	90078	42	69	5.0	0.5	Eighth
62430	119774	1	84	5.0	1.0	Quarter
62430	119774	1	72	5.0	1.0	Quarter

start_time	end_time	note
9182	90078	53
9182	33758	65
9182	62430	69
9182	202206	41
9182	62430	81
33758	62430	60
62430	90078	69
62430	119774	84
62430	119774	72
90078	119774	65
90078	119774	57
119774	145886	74

PRE-PROCESSING

Audio

- broken into fixed blocks
- builds a label matrix based on the block-level slices and the total possible MIDI note values for training (128)
- reformatted for use in the model

Labels

- MIDI note onset/offsets aligned with audio block timing
- converts .csv into a usable DataFrame format

CNN

Convolutional Neural Network

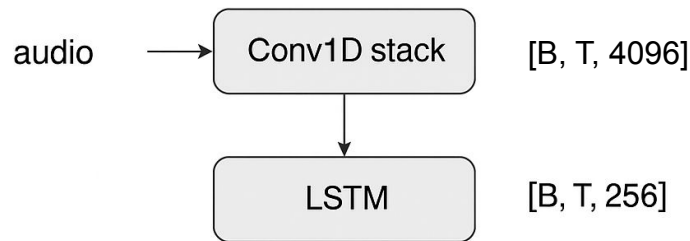
-Effective at recognizing local patterns and learning hierarchical features

-Generalizes well for wider numbers of pieces

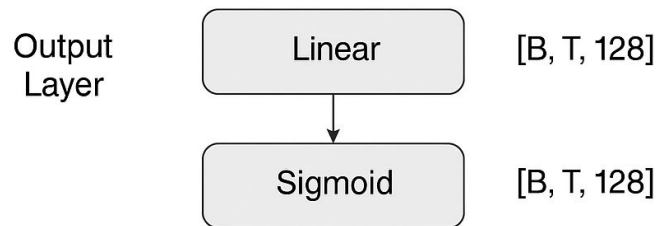
-research shows overfitting in Hidden Markov Model implementations

Waveform to MIDI Model

Convolutional Encoder



LSTM Sequence Model



MIDI note activations

DEFINED MODEL

Input ($B * T, 1, 4096$)



Convolution Encoder



LSTM Sequence Model



Fully Connected Layers



Softmax



Output

B: Batch size; T: Temporal length (audioblocks)
4096: block size (samples per block)

1D Convolutional layer: Extract time-domain features to each block

Long Short-Term Memory Sequence: Captures temporal dependencies across blocks

Maps temporal features to pitch activation space

Interpreted as a probability distribution over possible notes for each block

Generate output and convert to MIDI format

THANK YOU

R. Faran, Y. Wu, S. Xu

Automatic Polyphonic Music Transcription

Questions?

PRE-PROCESSING

Audio

- waveform broken into fixed-length blocks (10 blocks of 4096 samples each)
- builds a label matrix based on the block-level slices and the total possible MIDI note values for training (128)
 - [10 128]
- reformatted for use in the model

Labels

- MIDI note onset/offsets aligned with audio block timing
- converts .csv into a usable DataFrame format
- time converted into frame/block indices
- fills binary matrix [blocks, 128] with 1's where notes are active

Example Loaded Annotations of Training Label

Annotations DataFrame:

	start_time	end_time	instrument	note	start_beat	end_beat	\
0	9182	90078	43	53	4.0	1.5	
1	9182	33758	42	65	4.0	0.5	
2	9182	62430	1	69	4.0	1.0	
3	9182	202206	44	41	4.0	3.5	
4	9182	62430	1	81	4.0	1.0	
	note_value	duration					
0	Dotted Quarter	80896					
1	Eighth	24576					
2	Quarter	53248					
3	Whole	193024					
4	Quarter	53248					

CONVOLUTION ENCODER

Purpose:

- Each 4096-sample block is processed by a small CNN to extract local time-domain features.
- Operates on 1D audio waveform (like a short-time filter bank).

Steps:

- The Conv1D layers (with kernel size 9, stride 2) reduce the time resolution.
- MaxPool1D further reduces dimensionality.
- Final shape after convolution per block: [B, T, 256]

The CNN compresses each block (4096 samples) into a feature vector of length 256.

LSTM SEQUENCE MODEL

Purpose:

- Models temporal relationships between the blocks in a sequence.
- Each block has now been reduced to a 256-dim feature vector \rightarrow LSTM sees a sequence of vectors: $[B, T, 256]$.

Parameters:

- 2-layer LSTM
- Hidden size: 128
- Batch first: true (input shape is $[B, T, F]$)

Output:

- Produces output for each time step (block) \rightarrow shape $[B, T, 128]$

OUTPUT LAYER

Applies softmax across 128 MIDI pitches.

Interpreted as a **probability distribution over possible** notes for each block.

Each timestep gives a probability distribution for the 128 possible MIDI note values (127 notes or note off)

MASKED BINARY CROSS-ENTROPY LOSS

Computes BCE loss

-uses making padded positions to ensure only for active time steps

Averaged over valid tokens