

# High-resolution Piano Transcription with Pedals by Regressing Onset and Offset Times

Qiuqiang Kong, Bochen Li, Xuchen Song, Yuan Wan, Yuxuan Wang

**Abstract**—Automatic music transcription (AMT) is the task of transcribing audio recordings into symbolic representations. Recently, neural network-based methods have been applied to AMT, and have achieved state-of-the-art results. However, many previous systems only detect the onset and offset of notes frame-wise, so the transcription resolution is limited to the frame hop size. There is a lack of research on using different strategies to encode onset and offset targets for training. In addition, previous AMT systems are sensitive to the misaligned onset and offset labels of audio recordings. Furthermore, there are limited researches on sustain pedal transcription on large-scale datasets. In this article, we propose a high-resolution AMT system trained by regressing precise onset and offset times of piano notes. At inference, we propose an algorithm to analytically calculate the precise onset and offset times of piano notes and pedal events. We show that our AMT system is robust to the misaligned onset and offset labels compared to previous systems. Our proposed system achieves an onset F1 of 96.72% on the MAESTRO dataset, outperforming previous onsets and frames system of 94.80%. Our system achieves a pedal onset F1 score of 91.86%, which is the first benchmark result on the MAESTRO dataset. We have released the source code and checkpoints of our work at [https://github.com/bytedance/piano\\_transcription](https://github.com/bytedance/piano_transcription).

**Index Terms**—Piano transcription, pedal transcription, high-resolution.

## I. INTRODUCTION

Automatic music transcription (AMT) [1], [2], [3] is the task of transcribing audio recordings into symbolic representations [4], such as piano rolls, guitar fretboard charts and Musical Instrument Digital Interface (MIDI) files. AMT is an important topic of music information retrieval (MIR) and is a bridge between audio-based and symbolic-based music understanding. AMT systems have several applications, such as score following [5], audio to score alignment [6], and score-informed source separation [7]. In industry, AMT systems can be used to create music education software for music learners. For music production, AMT can be used to transcribe audio recordings into MIDI files for intelligent music editing. AMT systems can also be used for symbolic-based music information retrieval and can be used to analyze unarchived music, such as jazz improvisations.

Piano transcription is an essential task of AMT. The task is to transcribe piano solo recordings into music note events with pitch, onset, offset, and velocity. Piano transcription is a challenging task due to the high polyphony of music pieces. Early works of piano transcription [8], [9], [10] include using

discriminative models, such as support vector machines to predict the presence or absence of notes in audio frames [11]. To address the multiple pitch estimation problem, a probabilistic spectral smoothness principle was proposed in [12] for piano transcription. A combination of frequency domain and time domain method was proposed for piano transcription in [13], where authors assumed that signals are a linearly weighted sum of waveforms in a database of individual piano notes. Non-negative matrix factorizations (NMFs) and non-negative sparse codings [14] have been proposed to decompose spectrogram into polyphonic notes [15], where signals are decomposed into the multiplication of dictionaries and activations for transcribing polyphonic music. To model different onset and decay states of piano notes, an attack and decay system was proposed in [16]. Other AMT systems include using unsupervised learning method [17], connectionist approaches [18] and fast convolutional sparse coding methods [19].

Recently, neural networks have been applied to tackle the ATM problem. A deep belief network was proposed to learn feature representations for music transcription in [20]. Fully connected neural networks, convolutional neural networks (CNNs), and recurrent neural networks (RNNs) [21], [22], [23], [24] were proposed to learn regressions from audio input to labelled ground truths. Recently, onsets and frames systems [25], [26] were proposed to predict both onsets and frame-wise pitches of notes and have achieved state-of-the-art results in piano transcription. To improve the encoding of onset targets, several works including SoftLoc [27] and non-discrete annotations [28], [29], [30] have been proposed.

However, there are several limitations of previous AMT methods [25], [26]. First, those previous works used binarized values to encode onsets and offsets in each frame. Therefore, the transcription resolution in the time domain is limited to the frame hop size between adjacent frames. For example, a hop size of 32 ms was used in [25]. So that the transcription resolution is limited to 32 ms. In many scenarios, high-resolution transcription systems can be useful for music analysis. For example, Gobel [31], [32] analyzed the melody lead phenomenon in milliseconds. The second limitation is that the modeling of onsets and offsets can be improved. In [25], authors empirically used an onset length of 32 ms and claim that almost all onsets will end up spanning exactly two frames. However, there is a lack of explanation of labelling two frames as onsets performs better than other representations. The analysis from [16] shows that the attack of a piano note can last for several frames instead of only one frame. A note can be modeled in a more natural way with an attack, decay, sustain and release (ADSR) states [33]. In addition, the

Q. Kong, B. Li, X. Song, Y. Wan, Y. Wang are with ByteDance. (e-mail: kongqiuqiang@bytedance.com; bochenli@bytedance.com; xuchen.song@bytedance.com, wanyuan.0626@bytedance.com, wanyuxuan.11@bytedance.com). (Qiuqiang Kong is first and corresponding author.)

targets in [25] are sensitive to the misalignment between audio recordings and labels. For example, if onset is misaligned by one or several frames, then the training target [25] will be completely changed.

In this work, we propose a regression-based high-resolution piano transcription system that can achieve arbitrary resolution in the time domain for music transcription. In training, we propose regression-based targets to represent the time difference between the centre of a frame and its nearest onset or offset times. The physical explanation is that each frame is assigned a target of *how far* it is from its nearest onset or offset. Therefore, the information of precise onset or offset times have remained. At inference, we propose an analytical algorithm to calculate the precise onset or offset times with arbitrary time resolution. In evaluation, we investigate tolerances ranging from 2 ms to 100 ms for onset evaluation compared to the fixed tolerance of 50 ms [25]. We show that our proposed high-resolution piano transcription system achieves state-of-the-art results on the MAESTRO dataset [26]. We also show that our proposed regression-based targets are robust to the misaligned onsets and offsets. In addition, we develop a sustain pedal transcription system with our proposed method on the MAESTRO dataset, which has not been evaluated in previous works.

This paper is organized as follows. Section II introduces neural network based piano transcription systems. Section III introduces our proposed high-resolution piano transcription system. Section IV shows experimental results. Section V concludes this work.

## II. NEURAL NETWORK-BASED PIANO TRANSCRIPTION SYSTEMS

### A. Frame-wise Transcription Systems

Neural networks have been applied to tackle the piano transcription problem in previous works [21], [22], [23], [25]. First, audio recordings are transformed into log mel spectrograms as input features. We denote a log mel spectrogram as  $X \in \mathbb{R}^{T \times F}$ , where  $T$  is the number of frames, and  $F$  is the number of mel frequency bins. Then, neural networks, such as fully connected neural networks, CNNs, or RNNs are applied on the log mel spectrograms to predict the frame-wise presence probabilities of piano notes. Usually, a frame-wise roll  $I_{\text{fr}} \in \{0, 1\}^{T \times K}$  is used as a target for training [22], where  $K$  is the number of pitch classes, and is equivalent to 88 for piano transcription. The elements of  $I_{\text{fr}}$  have values of either 1 or 0, indicating the presence or absence of piano notes. Neural network-based methods [22] use a function  $f$  to map the log mel spectrogram of an audio clip to the frame-wise roll  $I_{\text{fr}}$ . We denote the neural network output as  $P_{\text{fr}} = f(X)$ , which has a shape of  $T \times K$ . The function  $f$  is modeled by a neural network with a set of learnable parameters. The following loss function is used to train the neural network [22]:

$$l_{\text{fr}} = \sum_{t=1}^T \sum_{k=1}^K l_{\text{bce}}(I_{\text{fr}}(t, k), P_{\text{fr}}(t, k)), \quad (1)$$

where  $l_{\text{fr}}$  is the frame-wise loss, and  $l_{\text{bce}}(\cdot, \cdot)$  is a binary cross-entropy function defined as:

$$l_{\text{bce}}(y, p) = -y \ln p - (1 - y) \ln(1 - p). \quad (2)$$

The target  $y \in \{0, 1\}$  is a binarized value and  $p \in [0, 1]$  is the predicted probability. A frame-wise piano transcription system is trained to predict the presence probability of notes in each frame [22].

At inference, we first calculate the log mel spectrogram of an audio recording. Then, the log mel spectrogram is input to the trained neural network to calculate  $f(X)$ . Finally, those predictions are post-processed to piano note events [22].

### B. Onsets and Frames Transcription System

For a piano, an onset is a hammer-string impact, which is equivalent to the beginning of a piano note event. An offset is a deactivated hammer-string impact [31], [32], for both MAESTRO dataset [26] and real-world recordings. One problem of the frame-wise transcription systems is that the transcribed frames need to be elaborately post-processed to piano notes [22]. In addition, those frame-wise piano transcription systems do not predict onsets explicitly, while the onsets carry rich information of piano notes. To address this problem, onsets and offsets dual objective system [25] was proposed to predict onsets and frames jointly. The onset and frame predictions are modeled by individual acoustic models containing several convolutional layers and long short term memory (LSTM) layers. The predicted onsets are used as conditional information to predict frame-wise outputs. We denote the predicted onset and frame outputs as  $P_{\text{on}}$  and  $P_{\text{fr}}$  respectively, where  $P_{\text{on}}$  and  $P_{\text{fr}}$  have shapes of  $T \times K$ . We denote the onset and frame targets as  $I_{\text{on}}$  and  $I_{\text{fr}}$  respectively, where  $I_{\text{on}}$  and  $I_{\text{fr}}$  also have shapes of  $T \times K$ . In [25], a joint frame and onset loss function was used to train the onsets and frames system:

$$l_{\text{note}} = l_{\text{on}} + l_{\text{fr}}, \quad (3)$$

where  $l_{\text{fr}}$  is the frame-wise loss defined in (1), and  $l_{\text{on}}$  is the onset loss defined as:

$$l_{\text{on}} = \sum_{t=1}^T \sum_{k=1}^K l_{\text{bce}}(I_{\text{on}}(t, k), P_{\text{on}}(t, k)). \quad (4)$$

One advantage of the onsets and frames system is that the onset predictions can be used as extra information to predict frame-wise outputs. The onsets and frames system has become a benchmark system for piano transcription.

## III. HIGH-RESOLUTION PIANO TRANSCRIPTION SYSTEM

Previous piano transcription systems introduced in Section II have several limitations. The methods in Section II predicts the presence or absence of onsets and frames in frames. Therefore, the transcription resolutions of those methods are limited to the hop size between adjacent frames. For example, the system [25] applies a hop size of 32 ms, so the transcription resolution in the time domain is limited to 32 ms. Second, for each piano note, previous systems [25] only label one or several frames of an onset or offset as 1, with other frames

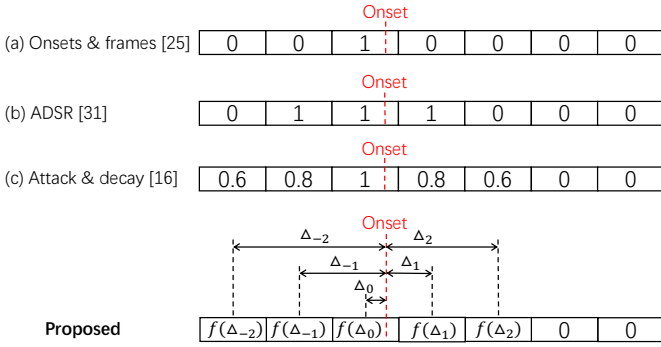


Fig. 1. Training targets of previous and our proposed piano transcription systems.

labelled as 0. The first row of Fig. 1 shows the onset targets used in [25]. The dashed vertical line shows the precise onset time of a note. The onsets and frames system [25] only assign one positive value to several consecutive frames indicating the onset of a piano note. However, this can be imprecise because it is unclear how many frames an attack will last. The onsets and frames system [25] labelled two consecutive frames as onsets empirically, which may not be optimal. In [33], ADSR states are used to model the onsets and offsets, and several neighbouring frames of an onset are labelled as 1, as shown in the second row of Fig. 1. The offset time of notes will not be changed by reverberation, but the waveform of audio recordings can be blurred, which will lead to the difficulty of offset detection.

In addition, the targets shown in the first row of Fig. 1 are sensitive to the misalignment of onset or offset labels. To explain, shifting the onset by one or several frames will lead to a completely different target. To mitigate this problem, Cheng et al. [16] proposed attack and decay targets to model the onset of piano notes shown in the third row of Fig. 1. Instead of only labeling the frames containing onsets as 1, the neighbouring frames of onsets are labelled with continuous values. Similar ideas have been proposed to tackle the pitch estimation [34] and music structure analysis [35] problems using smoothing filters for boundaries prediction.

Another problem of previous onset target representations [25], [33], [16] is that they do not reflect the precise onset or offset times of notes. The precise onset and offset times information is lost when quantizing onset and offset times into frames. We explain this in the first to the third rows of Fig. 1. The onset targets are unchanged when the precise onset times (dashed lines) are shifted within a frame. To achieve high-resolution piano transcription, the onset targets should be sensitive to the onset shifts in milliseconds. Furthermore, when the precise onset time is on the boundary between two frames, the target can be confusing. In addition, efforts to increase transcription resolution in the time domain by simply reducing frame hop size will take more computation cost, and the problem of limited transcription resolution still exists.

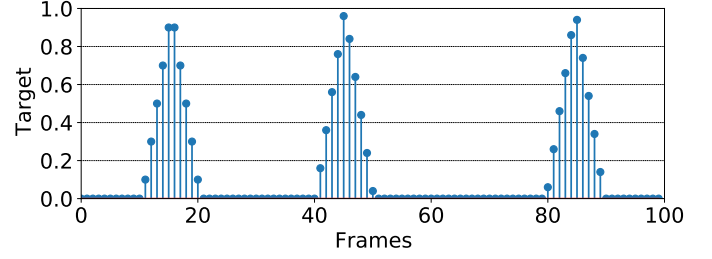


Fig. 2. High-resolution training targets of three notes with a same pitch.

### A. Regress Onset and Offset Times

We propose a high-resolution piano transcription system by predicting the continuous onset and offset times of piano notes instead of classifying the presence probabilities of onsets and offsets in each frame. This idea is inspired by the *You Look Only Once* (YOLO) [2] object detection method from computer vision. In YOLO, an image is split into grids. Then, each grid predicts a distance between the coordinate of the grid and the coordinate of an object. The distance to be predicted is a continuous value. Different from previous works [25], [33], [16], [34], [35], our proposed targets have a physical explanation that each frame is assigned a target of *how far* it is from its nearest onset or offset. The targets are calculated by the time distance between the centre of a frame and the precise onset time of a note. The bottom row of Fig. 1 shows the targets of our proposed high-resolution piano transcription system. We denote the frame hop size time as  $\Delta$ , and the time difference between the centre of a frame and its nearest onset time as  $\Delta_i$ , where  $i$  is the index of a frame. Negative  $i$  and positive  $i$  indicate previous and future frame indexes of an onset. Different from the targets of previous works [25], [33], [16] shown in the first to the third rows of Fig. 1, our proposed time difference  $\Delta_i$  contains precise onset and offset times information with arbitrary resolution. In training, we encode the time difference  $\Delta_i$  to targets  $g(\Delta_i)$  by a function  $g$ :

$$\begin{cases} g(\Delta_i) = 1 - \frac{|\Delta_i|}{J\Delta}, & |i| \leq J \\ g(\Delta_i) = 0, & |i| > J, \end{cases} \quad (5)$$

where  $J$  is a hyper-parameter controlling the sharpness of the targets. Larger  $J$  indicates “smoother” target, and smaller  $J$  indicates “sharper” target. Fig. 2 shows the visualization of onset targets of a pitch with  $J = 5$ . There are three piano notes in Fig. 2. Different from the attack and decay targets [16] shown in the third row of Fig. 1, the targets  $g(\Delta_i)$  in Fig. 2 contain precise onset times information of piano notes. For example, Fig. 2 shows that the onset of the first note is on the boundary of two adjacent frames, and the onsets of the second and third notes appear in different times. In training, both onset and offset regression targets are matrices with shapes of  $T \times K$ . We denote onset and offset regression targets as  $G_{\text{on}}$  and  $G_{\text{off}}$  respectively, to distinguish them from the binarized targets  $I_{\text{on}}$  and  $I_{\text{off}}$  in Section II. We denote the predicted onset and offset regression values as  $R_{\text{on}}$  and  $R_{\text{off}}$  respectively, to distinguish them from  $P_{\text{on}}$  and  $P_{\text{off}}$  in Section II. Both of regression based outputs  $R_{\text{on}}$  and  $R_{\text{off}}$  have values

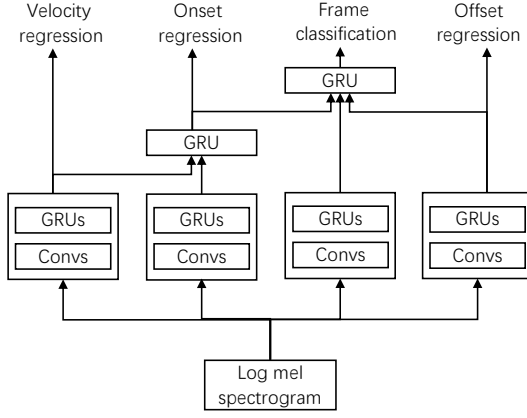


Fig. 3. High-resolution piano transcription system by regressing velocities, onsets, offsets and frames.

between 0 and 1. We define the onset regression loss  $l_{\text{on}}$  and offset regression loss  $l_{\text{off}}$  as:

$$l_{\text{on}} = \sum_{t=1}^T \sum_{k=1}^K l_{\text{bce}}(G_{\text{on}}(t, k), R_{\text{on}}(t, k)), \quad (6)$$

$$l_{\text{off}} = \sum_{t=1}^T \sum_{k=1}^K l_{\text{bce}}(G_{\text{off}}(t, k), R_{\text{off}}(t, k)). \quad (7)$$

Equation (6) and (7) use regression-based targets instead of classification-based targets in (4). To be consistent with the binary cross-entropy loss used in (1) and (4), we use binary cross-entropy in (6) and (7). Equation (6), and (7) are minimized when  $R_{\text{on}}$  equals  $G_{\text{on}}$  and  $R_{\text{off}}$  equals  $G_{\text{off}}$ .

### B. Velocity Estimation

Velocities of piano notes are correlated with the loudness of piano notes being played. Dannenberg and Goebel analyzed the correlation between velocity and loudness in [36] and [32]. In this work, we estimate the velocity of notes via estimating the MIDI velocities [26]. We build a velocity estimation submodule to estimate the velocities of transcribed notes. MIDI files represent velocities of notes using integers ranging from 0 to 127. Larger integers indicate loud notes and smaller integers indicate quiet notes. To begin with, we normalize the dynamic range of velocities from  $[0, 127]$  to  $[0, 1]$ . We denote the ground truth and predicted velocities as  $I_{\text{vel}}$  and  $P_{\text{vel}}$  respectively, where  $I_{\text{vel}}$  and  $P_{\text{vel}}$  have shapes of  $T \times K$ . Then, we define the velocity loss as:

$$l_{\text{vel}} = \sum_{t=1}^T \sum_{k=1}^K I_{\text{on}}(t, k) \cdot l_{\text{bce}}(I_{\text{vel}}(t, k), P_{\text{vel}}(t, k)). \quad (8)$$

Equation (8) shows that the ground truth onsets  $I_{\text{on}}(t, k)$  are used to modulate the velocity prediction. That is, we only predict velocities for onsets. One motivation is that the onset of piano notes carry rich information of their velocities, while the decay of piano notes carry less information of velocities than onsets. Similar to (6) and (7), binary cross-entropy is used to optimize (8). The loss  $l_{\text{vel}}$  is minimized when  $P_{\text{vel}}(t, k)$  equals  $I_{\text{vel}}(t, k)$ . At inference, we only predict velocities where onsets

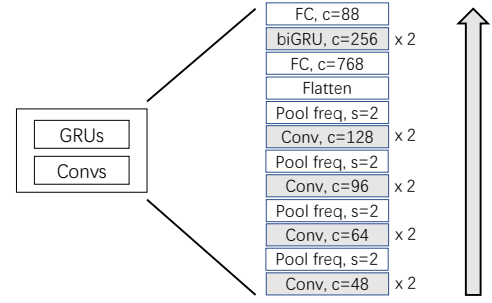


Fig. 4. Acoustic model. Pooling is only applied along the frequency axis. Letter “c” and “s” indicate the number of channels and the downsampling rate.

are detected. Finally, the predicted velocities are scaled from  $[0, 1]$  back to integers of  $[0, 127]$ .

### C. Entire System

Fig. 3 shows the framework of our proposed high-resolution piano transcription system. First, an audio clip is transformed into a log mel spectrogram with a shape of  $T \times F$  as input feature [25], where  $F$  is the number of mel frequency bins. There are four submodules in Fig. 3 from left to right: a velocity regression submodule, an onset regression submodule, a frame-wise classification submodule, and an offset regression submodule. Each submodule is modeled by an acoustic model [25]. In our system, we model each acoustic model with several convolutional layers followed by bidirectional gated recurrent units (biGRU) layers. The convolutional layers are used to extract high-level information from the log mel spectrogram, and the biGRU layers are used to summarize long-time information of the log mel spectrogram. Then, a time-distributed fully connected layer is applied after the biGRU layer to predict regression results. The outputs of all acoustic models have dimensions of  $T \times K$ . Fig. 4 shows the neural network layers of the acoustic models used in Fig. 3. We will describe the detailed configuration of the acoustic models in Section IV-C.

Fig. 3 shows that the predicted velocities are used as conditional information to predict onsets. One motivation is that the velocity of a piano note can be helpful to detect its corresponding onset. For example, if the velocity of a note is low, then the system will attend more to the frame to detect the onset. This simulates human beings who will listen carefully to the notes with low velocity. We concatenate the outputs of the velocity regression submodule and the onset regression submodule along the frequency dimension and use this concatenation as an input to a biGRU layer to calculate the final onset predictions. Similarly, we concatenate the outputs of the onset regression and offset regression submodules and use this concatenation as the input to a biGRU layer to calculate the frame-wise predictions. The total loss function to train our proposed piano transcription system consists of four parts:

$$l_{\text{note}} = l_{\text{fr}} + l_{\text{on}} + l_{\text{off}} + l_{\text{vel}}, \quad (9)$$

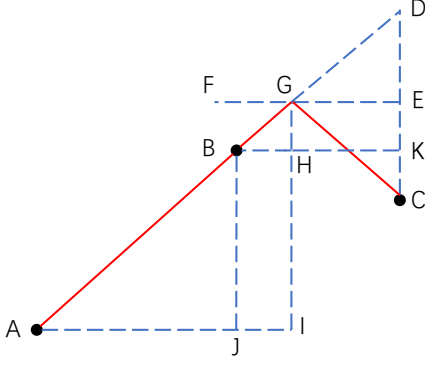


Fig. 5. Demonstration of calculating the precise onset or offset time of a note. The points  $A$ ,  $B$  and  $C$  are the predictions of three frames. The point  $G$  is the calculated precise onset or offset time.

where  $l_{fr}$ ,  $l_{on}$ ,  $l_{off}$  and  $l_{velocity}$  are described in (1), (6), (7) and (8) respectively. We simply weight all those losses equally, which works well in our experiment.

#### D. Inference

At inference, we input the log mel spectrogram of an audio recording into the trained piano transcription system to calculate the frame-wise prediction, onset regression, offset regression, and velocity regression outputs. Then, we propose an algorithm to process those outputs to high-resolution note events, where each note event can be represented by a quadruple of  $\langle$ piano note, onset time, offset time, velocity $\rangle$ .

Fig. 5 shows the strategy of calculating precise onset or offset times of piano notes. The horizontal coordinate of  $A$ ,  $B$ ,  $C$  are the centre times of three adjacent frames, where  $B$  is the frame with a local maximum onset prediction value. Previous works [25] regard the time of  $B$  as the onset time. In this case, the transcription resolution is limited to the hop size between adjacent frames. Ideally, the precise onset time should be  $G$ . We propose to calculate the precise onset time of  $G$  as follows. First, we detect  $A$ ,  $B$ ,  $C$  where  $B$  is a local maximum frame. If the vertical value of  $B$  is larger than an onset threshold, then we say there exists an onset near the frame  $B$ . Next, we analytically calculate the precise time of  $G$  where  $G$  satisfies  $AG$  and  $CG$  are symmetric along the vertical line  $GI$ .

Without loss of generality, we assume the output value of  $C$  is larger than  $A$ . We denote the coordinate of  $A$ ,  $B$  and  $C$  as  $(x_A, y_A)$ ,  $(x_B, y_B)$  and  $(x_C, y_C)$  respectively. We show that the time difference between  $B$  and  $G$  is:

$$BH = \frac{x_B - x_A}{2} \frac{y_C - y_A}{y_B - y_A}. \quad (10)$$

*Proof.* Extend  $AB$  to  $D$  where  $CD$  is a vertical line. Take the median point of  $CD$  as  $E$ . Draw a horizontal line  $EF$  cross  $AD$  at  $G$ . Then, we calculate  $BH$ . We have  $\triangle BGH \sim \triangle ABJ$ , so  $BH = \frac{AJ \cdot GH}{BJ}$ . We know that  $AJ = x_B - x_A$ ,  $BJ = y_B - y_A$  and  $GH = DK - DE$ , and we know  $AJ = BK$ , so  $\triangle ABJ \cong \triangle BDK$ , so  $DK = y_B - y_A$ . Then, we have  $DE = \frac{CD}{2} = \frac{DK + CK}{2} = y_B - \frac{y_A + y_C}{2}$ , so  $GH = \frac{y_C - y_A}{2}$ , so  $BH = \frac{x_B - x_A}{2} \frac{y_C - y_A}{y_B - y_A}$ .  $\square$

---

#### Algorithm 1 Onset and offset times detection.

---

```

1: Inputs:  $R_{on}(t, k)$ ,  $R_{off}(t, k)$ ,  $P_{fr}(t)$ ,  $P_{vel}(t, k)$ ,  $\theta_{on}$ ,  $\theta_{off}$  and  $\theta_{fr}$ .
2: Outputs: Detected onset and offset times.
3: for  $k = 1, \dots, K$  do
4:   for  $t = 1, \dots, T$  do
5:     # Detect note onset.
6:     if  $R_{on}(t, k) > \theta_{on}$  and  $R_{on}(t, k)$  is local maximum then
7:       Note onset of pitch  $k$  is detected. The precise onset time is refined by (10) or (11).
8:       Calculate the velocity of the note by  $P_{vel}(t, k) \times 128$ .
9:     end if
10:    # Detect note offset.
11:    if  $(R_{off}(t, k) > \theta_{off}$  and  $R_{off}(t, k)$  is local maximum) or  $P_{fr}(t, k) < \theta_{fr}$  then
12:      Note offset is  $k$  detected. The precise offset time is refined by (10) or (11).
13:    end if
14:  end for
15: end for

```

---

In another case, if the output value of  $A$  is larger than  $C$ , then:

$$BH = \frac{x_C - x_B}{2} \frac{y_A - y_C}{y_B - y_C}. \quad (11)$$

By this means, we can calculate the precise onset or offset times of piano notes. We describe onset and offset times detection in Algorithm 1. A frame is detected to contain an onset if they are over an onset threshold  $\theta_{on}$  and is a local maximum. Then, the precise onset time is calculated by (10) or (11). The velocity of an onset is obtained by scaling the predicted velocity from  $[0, 1]$  to a range of  $[0, 127]$ . For each detected onset, an offset is detected if the offset regression output is over an offset threshold  $\theta_{off}$  or any frame prediction outputs are lower than a frame threshold  $\theta_{fr}$ . When consecutive onsets of the same pitch are detected, their previous onsets are truncated by adding offsets. Our proposed method can not detect consecutive notes shorter than 4 frames (40 ms) because of the local maximum detection algorithm. Still, repeating a note shorter than 40 ms rarely happens in real piano performance. Finally, all onsets and offsets are paired to constitute piano notes.

#### E. Sustain Pedal Transcription

In this section, we will show our proposed regression-based transcription system can be applied to sustain pedal detection. The sustain pedal transcription system is similar to the piano note transcription system. The only difference is that the sustain pedal transcription system has output shapes of  $T \times 1$  instead of  $T \times 88$ , where the outputs indicate the presence probability of sustain pedals. Sustain pedals are important for piano performance. When pressed, the sustain pedal sustains all damped strings on a piano by moving all dampers away from the strings and allows strings to vibrate freely. All notes being played will continue to sound until the pedal is

released. However, many previous piano transcription systems [25], [16], [33] did not incorporate sustain pedal transcription. On the other hand, some sustain pedal transcription systems [37] do not include piano notes transcription. In [37], a convolutional neural network was used to detect piano pedals in each frame. Still, there is a lack of benchmark sustain pedal transcription systems on the MAESTRO dataset [26].

In this section, we propose a sustain pedal transcription system with our proposed high-resolution transcription system. The sustain pedal transcription system is separate from the note transcription system and is trained separately. Training the note and sustain pedal systems separately leads to better transcription performance and also has the advantage of reducing memory usage. After training the note transcription and sustain pedal transcription systems, we combine them into a unified model for release. In the MIDI format, sustain pedals are represented with integer values ranging from 0 to 128. MIDI values larger than 64 are regarded as “on” and MIDI values smaller than 64 are regarded as “off”.

To simplify the sustain pedal transcription problem, we only classify the “on” and “off” states of sustain pedals and do not consider advanced sustain pedal techniques such as half pedals. We denote the pedal onset regression target, offset regression target, and frame-wise target as  $G_{\text{ped\_on}} \in [0, 1]^T$ ,  $G_{\text{ped\_off}} \in [0, 1]^T$ , and  $I_{\text{ped\_fr}} \in \{0, 1\}^T$ , respectively. The onset and offset regression targets  $G_{\text{ped\_on}}$  and  $G_{\text{ped\_off}}$  are obtained by (5), and have continuous values between 0 and 1. The frame-wise targets have binarized values between 0 and 1. We apply acoustic models described in Section III-C to predict the onset regression, offset regression and frame-wise outputs of pedals. We denote the predicted onsets, offsets, and frame-wise values as  $R_{\text{ped\_on}}(t)$ ,  $R_{\text{ped\_off}}(t)$  and  $P_{\text{ped\_fr}}(t)$ , respectively. We use the following loss function to train the sustain pedal transcription system:

$$l_{\text{ped\_on}} = \sum_{t=1}^T l_{\text{bce}}(G_{\text{ped\_on}}(t), R_{\text{ped\_on}}(t)), \quad (12)$$

$$l_{\text{ped\_off}} = \sum_{t=1}^T l_{\text{bce}}(G_{\text{ped\_off}}(t), R_{\text{ped\_off}}(t)), \quad (13)$$

$$l_{\text{ped\_fr}} = \sum_{t=1}^T l_{\text{bce}}(I_{\text{ped\_fr}}(t), P_{\text{ped\_fr}}(t)). \quad (14)$$

Then, the total loss function is calculated by:

$$l_{\text{ped}} = l_{\text{ped\_fr}} + l_{\text{ped\_on}} + l_{\text{ped\_off}}. \quad (15)$$

At inference, we propose Algorithm 2 to process the sustain pedal prediction outputs into sustain pedal events. Sustain pedal onsets detection is different from notes onset detection. The press of a sustain pedal is usually before the press of piano notes so  $R_{\text{ped\_on}}(t)$  can be difficult to detect. We design a sustain pedal detection system so that a pedal onset is detected when the frame-wise prediction  $R_{\text{ped\_fr}}(t)$  is over a threshold  $\theta_{\text{ped\_on}}$ . A pedal offset is detected if the pedal offset prediction  $R_{\text{ped\_off}}(t)$  is higher than an offset threshold  $\theta_{\text{ped\_off}}$  or the frame-wise prediction  $P_{\text{ped\_fr}}(t)$  is lower than a threshold  $\theta_{\text{ped\_fr}}$ .

---

**Algorithm 2** Sustain pedal onset and offset times detection.

---

```

1: Inputs:  $R_{\text{off}}(t)$ ,  $R_{\text{fr}}(t)$ ,  $\theta_{\text{ped\_off}}$ ,  $\theta_{\text{ped\_fr}}$ .
2: Outputs: Detected pedal onset and offset times.
3: for  $t = 1, \dots, T$  do
4:   # Detect pedal onset.
5:   if  $R_{\text{ped\_fr}}(t) > \theta_{\text{ped\_on}}$  and  $R_{\text{ped\_fr}}(t) > R_{\text{ped\_fr}}(t - 1)$ 
     then
6:     Pedal onset is detected.
7:   end if
8:   # Detect pedal offset.
9:   if  $R_{\text{ped\_off}}(t) > \theta_{\text{ped\_off}}$  or  $P_{\text{ped\_fr}}(t) < \theta_{\text{ped\_fr}}$  then
10:    Pedal offset is detected. The precise offset time is
       refined by (10) or (11).
11:   end if
12: end for

```

---

## IV. EXPERIMENTS

### A. Dataset

To compare with previous piano transcription systems, we use the MAESTRO dataset V2.0.0 [26], a large-scale dataset containing paired audio recording and MIDI files to train and evaluate our proposed piano transcription system. The MAESTRO dataset contains piano recordings from the International Piano-e-Competition. Pianists performed on Yamaha Disklaviers concert-quality acoustic grand pianos integrated with high-precision MIDI capture and playback system. The MAESTRO dataset contains over 200 hours of solo piano recordings. Those audio recordings and MIDI files are aligned with a time resolution of around 3 ms introduced by [26]. Each music recording contains meta-information, including the composer, title, and year of the performance. MAESTRO dataset consists of training, validation and testing subsets.

### B. Preprocessing

We use Python and PyTorch deep learning toolkit [38] to develop our systems. All stereo audio recordings are converted into mono and are resampled to 16 kHz following [26]. The cutoff frequency of 16 kHz covers the frequency of the highest note C<sub>8</sub> on a piano of 4186 Hz. We split audio recordings into 10-second clips. Then, a short-time Fourier transform with a Hann window size 2048 is used to extract the spectrogram. Mel banks with 229 banks and cutoff frequencies between 30 Hz and 8000 Hz are used to extract the log mel spectrogram [26]. We use a hop size of 10 ms between frames. For a 10-second audio clip, an input log mel spectrogram has a shape of  $1001 \times 229$ , where the extra one frame comes from that audio clips are padded with half windows in both sides of a signal before feature extraction. Log mel spectrograms are extracted on the fly using the TorchLibrosa toolkit [39]. For general experiments, we set the hyper-parameter  $J = 5$ . That is, each onset or offset will affect the regression values of  $2 \times J = 10$  frames. We also investigate piano transcription systems with different  $J$  in our experiments. The first row of Fig. 6 shows an example of the log mel spectrogram of a 5-second audio clip. The second and third rows show the frame-wise targets and frame-wise predictions of the audio clip. The fourth and fifth



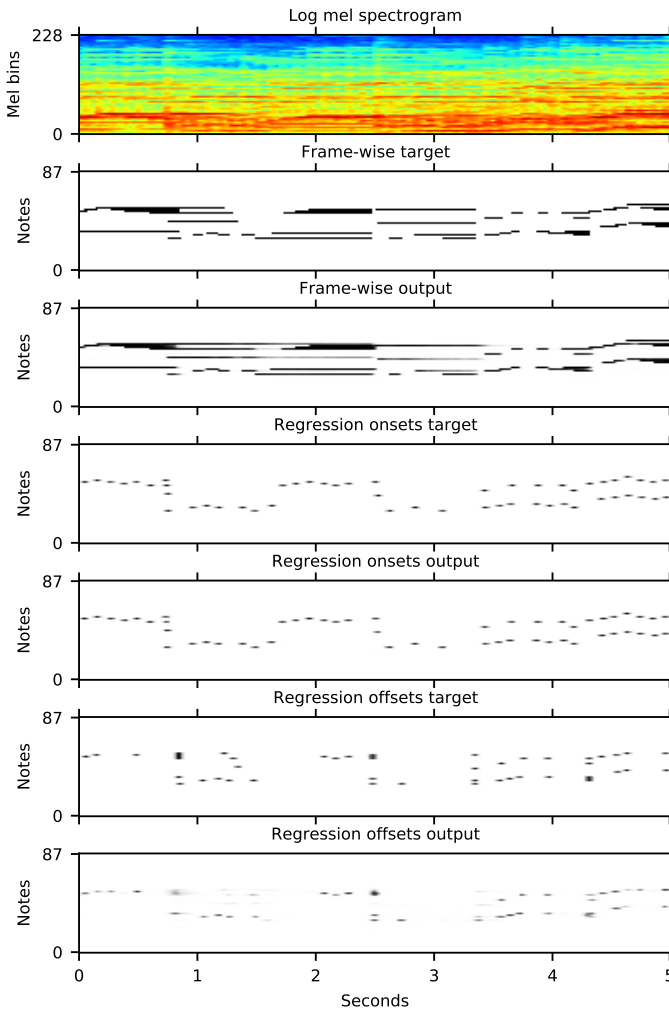


Fig. 6. From top to bottom: Log mel spectrogram of a 5-second audio clip; frame-wise targets; frame-wise outputs; onset regression targets; onset regression outputs; offset regression targets; offset regression outputs. The music segment is: Prelude and Fugue in D Major, WTC I, BWV 850, J. S. Bach, 2'35.

rows show the regression onset targets and onset predictions of the audio clip. The sixth and seventh rows show the regression offset targets and offset predictions of the audio clip.

### C. Model Architecture

After extracting the log mel spectrogram features, we apply a batch normalization [41] layer immediately on the individual frequency bins of the log mel spectrogram [39] to standardize the input. Then, acoustic models shown in Fig. 3 modeled by CRNNs are applied to predict the velocity regression, onset regression, frame-wise classification, and offset regression outputs. All acoustic models have the same architecture, where each acoustic model consists of four convolutional blocks and two bidirectional biGRU layers. Each convolutional block consists of two convolutional layers with kernel sizes  $3 \times 3$ . Batch normalization [41] and ReLU nonlinearity [42] are applied after each linear convolutional operation to stabilize training and to increase the nonlinear representation ability of the system. The four convolutional blocks have output feature

map numbers of 48, 64, 92 and 128 respectively. After each convolutional block, feature maps are averagedly pooled by a factor of 2 along the frequency axis to reduce the feature map sizes. We do not apply pooling along the time axis to retain the transcription resolution in the time domain.

After convolutional layers, feature maps are flattened along the frequency and channel axes and are input to a fully connected layer with 768 output units. Then, two biGRU layers with hidden sizes 256 are applied, followed by an additional fully connected layer with 88 sigmoid outputs. Dropout [43] with rates of 0.2 and 0.5 are applied after convolutional blocks and fully connected layers to prevent the systems from overfitting. Fig. 3 shows that the velocity and onset regression outputs are concatenated and are input to a biGRU layer. The biGRU layer contains 256 hidden units and is followed by a fully-connected layer with 88 sigmoid outputs to predict the final regression onsets. Similarly, the onset regression, offset regression, and frame-wise classification outputs are concatenated and are input to a biGRU layer. The biGRU contains 256 hidden outputs and is followed by a fully-connected layer with 88 sigmoid outputs to predict the final frame-wise output. The note transcription system consists of 20,218,778 trainable parameters. The training of the note transcription system applies a loss function described in (9). The sustain pedal transcription submodule has the same acoustic model architecture as the note transcription submodule, except there is only one output instead of 88 outputs. The pedal onset regression, offset regression, and frame-wise classification are modeled by individual acoustic models. The training of the sustain pedal transcription system applies a loss function described in (15).

We use a batch size 12, and an Adam [44] optimizer with a learning rate of 0.0005 for training. The learning rate is reduced by a factor of 0.9 every 10 k iterations in training. Systems are trained for 200 k iterations. The training takes four days on a single Tesla-V100-PCIE-32GB GPU card. At inference, we set onset, offset, frame-wise, and pedal thresholds to 0.3. All hyper-parameters are tuned on the validation set. The outputs are post-processed to MIDI events described in III-D.

### D. Evaluation

We evaluate our proposed piano transcription system on the test set of the MAESTRO dataset. We compare our system with previous onsets and frames system [25] and the adversarial onsets and frames system [40]. The system [25] is an improvement to the onsets and frames system [25] and also used quantized targets for training. For a fair comparison with our proposed system with previous systems, we re-implemented the onsets and frames system trained with hard labels of 0 and 1. The results are shown in the third row of Table I. The numbers of our re-implemented system are slightly different from [25] due to different data augmentation strategies, different data pre-processing, data post-processing strategies, and deep learning toolkits. There are four types of evaluation metrics for piano transcription evaluation, including frame-wise evaluation, note evaluation

TABLE I  
TRANSCRIPTION RESULTS EVALUATED ON THE MAESTRO DATASET

	FRAME			NOTE			NOTE W/ OFFSET			NOTE W/ OFFSET & VEL.		
	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)
Onsets & frames [26]	93.10	85.76	89.19	97.42	92.37	94.80	81.84	77.66	79.67	78.11	74.13	76.04
Onsets & frames [reproduced]	86.63	90.89	88.63	99.52	89.23	93.92	80.43	72.27	75.99	79.51	71.48	75.14
Adversarial onsets & frames [40]	93.10	89.80	<b>91.40</b>	98.10	93.20	95.60	83.50	79.30	81.30	82.30	78.20	80.20
Regress cond on.	86.94	90.15	88.42	98.43	94.84	96.57	80.00	77.08	78.50	78.64	75.79	77.17
Regress cond on. & off.	88.91	90.28	89.51	98.53	94.81	96.61	83.81	80.70	82.20	82.36	79.33	80.79
Regress cond on. & off. & vel.	88.71	90.73	89.62	98.17	95.35	<b>96.72</b>	83.68	81.32	<b>82.47</b>	82.10	79.80	<b>80.92</b>

TABLE II  
TRANSCRIPTION RESULTS EVALUATED WITH DIFFERENT HYPER-PARAMETER  $J$

	FRAME			NOTE			NOTE W/ OFFSET			NOTE W/ OFFSET & VEL.		
	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)
Hyper-parameter $J = 2$	85.36	91.76	88.35	98.94	93.68	96.19	81.42	77.15	79.19	80.02	75.85	77.84
Hyper-parameter $J = 5$	87.62	90.88	<b>89.14</b>	98.15	95.15	<b>96.61</b>	82.92	80.42	<b>81.63</b>	81.35	78.92	<b>80.10</b>
Hyper-parameter $J = 10$	86.89	90.89	88.75	97.60	95.39	96.47	81.48	79.66	80.55	79.82	78.05	78.92
Hyper-parameter $J = 20$	86.56	90.41	88.34	96.23	95.15	95.67	77.79	76.95	77.36	76.24	75.43	75.82

with onset, note evaluation with both onset and offset, and note evaluation with onset, offset, and velocity. Following [25], a tolerance of 50 ms is used for onset evaluation. A tolerance of 50 ms and an offset ratio of 0.2 are used for offset evaluation. A velocity tolerance of 0.1 is used for velocity evaluation, which indicates that estimated notes are considered correct if, after scaling and normalization velocities to a range of 0 to 1, they are within the velocity tolerance of a matched reference note.

The first to the third rows of Table I show the results of the onsets and offsets system [26], the reproduced onsets and offsets system, and the adversarial system [40] for a fair comparison. The fourth row shows that our proposed regression-based system only using onset as the condition improves the onsets and frames system [26] note F1 score from 94.80% to 96.57%. The fifth row shows that using both onset and offset as conditions achieves an onset F1 score of 96.61%. The sixth row shows that using onset and offset to condition the frame prediction and using velocity to condition the onset prediction further improves the note F1 score to 96.72%. We also evaluated the system performance of different runs and observed that the standard variance of frame F1 scores, onset F1 scores, onset F1 scores evaluated with offsets, and onset F1 scores evaluated with offsets and velocities are  $\pm 0.04\%$ ,  $0.03\%$ ,  $0.08\%$ , and  $0.10\%$ . Our proposed high-resolution system improves the note F1 score evaluated with offsets from 79.67% to 82.47% and improves the note F1 score evaluated with both offset and velocity from 76.04% to 80.92%. The first row of Fig. 6 shows the log mel spectrogram of an audio clip. The second and third rows show the frame-wise target and frame-wise system output. The fourth and fifth rows show the regression onset target and regression onset output. The sixth and seventh row show the regression offset target and regression offset output. The onset and offset regression targets are calculated by (5). Fig. 6 shows that our

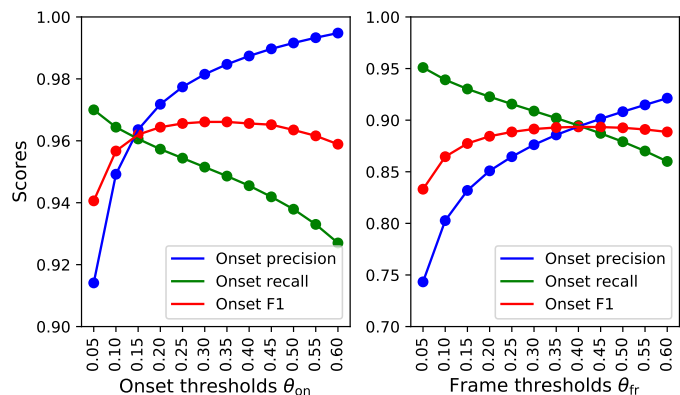


Fig. 7. Left: Onset-level precision, recall and F1 score evaluated with different onset thresholds  $\theta_{on}$ . Right: Frame-level precision, recall and F1 score evaluated with different frame thresholds  $\theta_{fr}$ .

proposed system performs well on transcribing a 5-second audio clip.

To evaluate different thresholds for piano transcription, we experiment with onset thresholds and frame thresholds from 0.05 to 0.60. Fig. 7 shows the precision, recall, and F1 score of the piano note transcription system evaluated with different onset thresholds and offset thresholds. Fig. 7 shows that F1 scores are similar when onset thresholds and frames thresholds are between 0.2 and 0.4. Higher thresholds lead to higher precision but lower recall. According to those observations, we set all thresholds to 0.3 in our work. Table II shows the piano transcription results with different hyper-parameters  $J$ . The evaluation scores of systems with different  $J$  are similar, indicating that  $J$  does not need to be tuned elaborately in our proposed system. Overall, setting  $J$  to 5 slightly outperforms other configurations.

To show that our proposed regression-based piano transcription system is robust to the misalignment of labels, we



TABLE III  
TRANSCRIPTION RESULTS EVALUATED WITH MISALIGNED LABELS

	FRAME			NOTE			NOTE W/ OFFSET			NOTE W/ OFFSET & VEL.		
	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)
Onsets & frames (misaligned labels)	80.93	90.93	85.54	65.59	93.06	76.52	44.40	63.36	51.92	40.41	57.64	47.25
Regress cond on. & off. & vel. (misaligned labels)	84.65	91.36	<b>87.79</b>	98.65	94.30	<b>96.39</b>	80.59	77.09	<b>78.77</b>	77.35	74.02	<b>75.62</b>

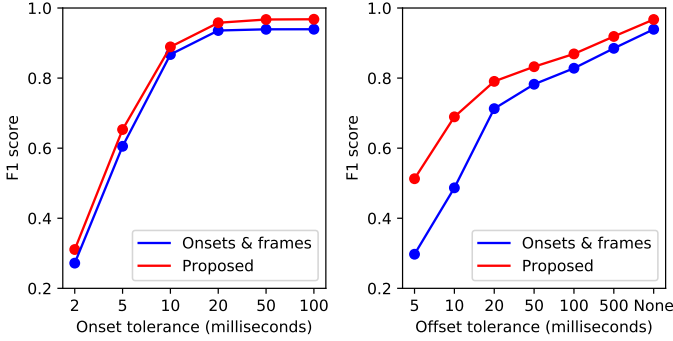


Fig. 8. Left: F1 score evaluated with various onset tolerances. Right: F1 score evaluated with various offset tolerances.

randomly shift the labels of onsets and offsets with a uniform distribution between  $-A$  ms and  $+A$  ms. To show that our system can generalize well to misaligned labels, we set  $A$  to 50 ms which is the default onset tolerance for onset evaluation used by the mir\_eval toolkit [45]. The seventh row of Table III shows that the note F1 score of [25] decreases from 93.92% to 76.52% when trained with misaligned labels. One explanation is that, the system [25] is sensitive to misaligned labels. On the other hand, the second row of Table III shows that our proposed regression-based system achieves a note F1 of 96.39%, compared to the system trained with correct labels of 96.72%. Our system achieves an F1 of 75.62% when evaluated with offsets and velocities, compared to the system trained with correct labels of 80.92%. Those results indicate that our proposed system is robust to misaligned labels.

To mathematically demonstrate that our proposed regression targets are robust to misaligned labels, we denote the distribution of misaligned labels as  $q(t)$  and denote the target of a note as  $f(t)$ , where  $f(t)$  can be either an onset target [25] shown in the top left of Fig. 9 or our proposed regression target (5) shown in the top right of Fig. 9. When using the binary cross-entropy loss function (2) for training, the optimal estimation of targets is  $u(t) = f(t) * q(t)$  where the symbol  $*$  is a convolution operation. In absence of misaligned label, there is  $q(t) = \delta(t)$  and the optimal estimation  $u(t)$  is equivalent to  $f(t)$ . When  $q(t) \sim [-A, A]$  is a uniform distribution, the optimal estimation  $u(t)$  of [25] and our proposed regression-based method are shown in the bottom row of Fig. 9. The bottom left of Fig. 9 shows that the precise onset times are hard to obtain when using the targets of [25] in training. In contrast, the bottom right of Fig. 9 shows that the precise onset time can be obtained when using the regression-based method in training.

We evaluate piano transcription performance with different

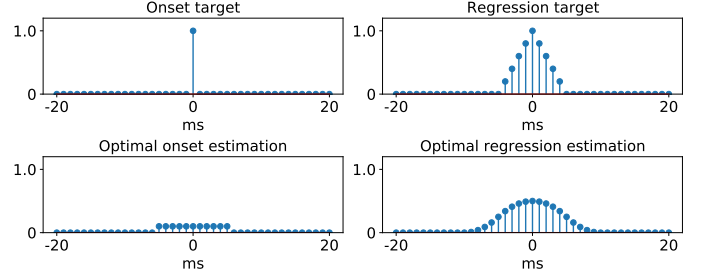


Fig. 9. Top left: onset target of [25]; Top right: our proposed regression target; Bottom left: optimal onset estimation with misaligned labels; Bottom right: optimal regression estimation of misaligned labels.

TABLE IV  
ONSETS EVALUATION WITH DIFFERENT ONSET TOLERANCES

	ONSETS & FRAMES			PROPOSED		
	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)
100 ms	99.54	89.24	93.94	98.24	95.42	<b>96.79</b>
50 ms	99.52	89.23	93.92	89.17	95.35	<b>96.72</b>
20 ms	99.14	88.92	93.58	97.22	94.45	<b>95.79</b>
10 ms	91.74	92.53	86.73	90.19	87.69	<b>88.91</b>
5 ms	63.89	57.69	60.53	66.22	64.47	<b>65.32</b>
2 ms	28.69	25.93	27.19	31.49	30.68	<b>31.08</b>

onset and offset tolerances that have not been evaluated in previous works [25], [26], [40]. The tolerances to be estimated range from 10 ms to 500 ms. Table IV and the left part of Fig. 8 show that with an onset tolerance of 2 ms, our system achieves an onset F1 score of 31.08%. The F1 score increases to 88.91% when onset tolerance is 10 ms, and increases to 96.79% when onset tolerance is 100 ms. The F1 scores of our proposed system outperform the onsets and offsets system [25] in all onset tolerances. Table V and the right part of Fig. 8 show the note F1 score evaluated with offset tolerances

TABLE V  
ONSETS AND OFFSETS EVALUATION WITH DIFFERENT OFFSET TOLERANCES

	ONSETS & FRAMES			REGRESS & ON. & OFF.		
	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)
500 ms	93.71	84.13	88.49	93.25	90.59	<b>91.88</b>
200 ms	87.65	78.75	82.81	88.19	85.68	<b>86.90</b>
100 ms	82.80	74.41	78.24	84.48	82.08	<b>83.25</b>
50 ms	75.47	67.76	71.28	80.24	77.95	<b>79.06</b>
20 ms	51.66	46.19	48.68	69.96	67.97	<b>68.94</b>
10 ms	31.64	28.14	29.73	52.04	50.57	<b>51.28</b>

TABLE VI  
PEDAL TRANSCRIPTION EVALUATED ON THE TEST SET OF MAESTRO DATASET.

	FRAME			EVENT			EVENT W/ OFFSET		
	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)
Liang	74.29	90.01	79.12	-	-	-	-	-	-
Onsets & frames [our-implemented]	94.30	94.42	94.25	93.20	90.26	91.57	86.94	84.28	85.47
Proposed	94.30	94.42	94.25	91.59	92.41	<b>91.86</b>	86.36	87.02	<b>86.58</b>
Onsets & frames (misaligned labels)	93.62	94.14	93.77	92.71	85.48	88.69	83.17	77.03	79.78
proposed (misaligned labels)	94.41	93.29	93.73	91.62	91.17	<b>91.23</b>	86.33	85.83	<b>85.94</b>

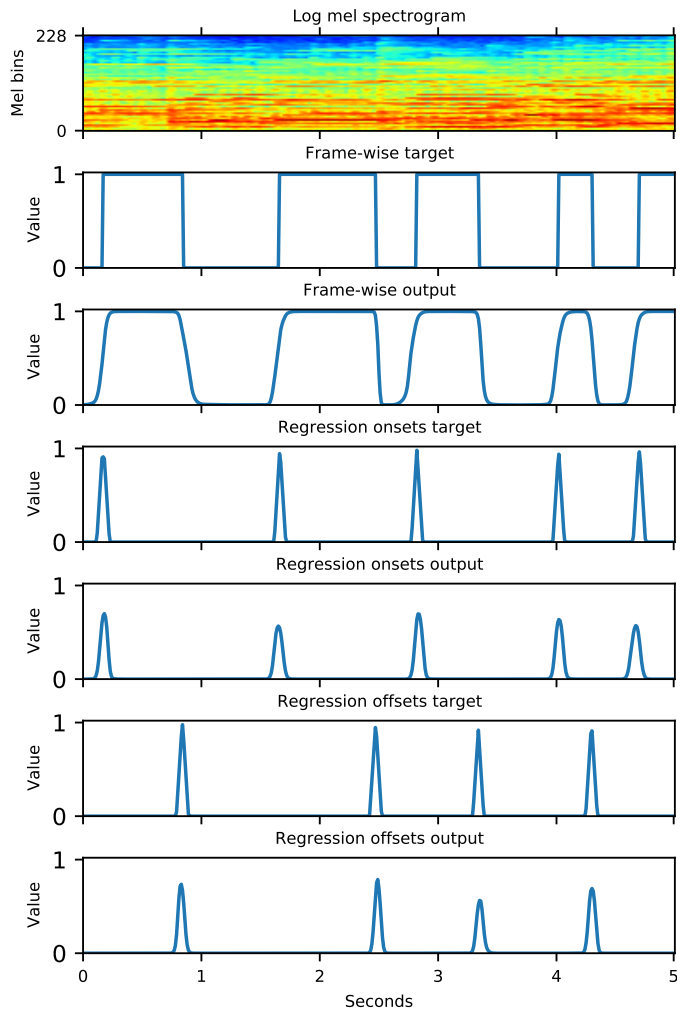


Fig. 10. Log mel spectrogram of a 5-second audio clip; pedal onset target; pedal onset output; pedal offset target; pedal offset output; pedal frame target; pedal frame output.

ranging from 10 ms to 500 ms when fixing onset tolerance to 50 ms. The *None* label in Fig. 8 equals to evaluation without offsets. Our system achieves a note F1 score of 51.28% with an offset tolerance of 10 ms. The F1 score increases to 91.88% with a tolerance of 500 ms. Our proposed system outperforms the onsets and frames system [25] in all offset tolerances. The experiments show that our proposed system can achieve higher transcription resolution than [25], [26].

We evaluate the pedal transcription results as follows.

Previous onsets and frames system [25] does not include sustain pedal transcription. Therefore, we implemented a pedal transcription with the onsets and frames system that has the same model architecture as our regression-based system for a fair comparison. In testing, for piano pieces without sustain pedals, the sustain pedal piano transcription system will produce no pedal events. We implemented the CNN-based method in [37], and achieves a frame-wise precision, recall, and F1 score of 74.29%, 90.01%, and 79.12%, respectively. Table VI shows the pedal transcription result. Our system achieves an event-based F1 of 91.86% evaluated with a pedal onset tolerance of 50 ms and achieves an event-based F1 of 86.58% evaluated with both onset and offset tolerances of 50 ms and offset ratio of 0.2, outperforming our implemented onsets and frames system of 91.57% and 85.47% respectively. As far as we know, we are the first to evaluate sustain pedal transcription on the MAESTRO dataset. The first row of Fig. 10 shows the log mel spectrogram of the audio clip that is the same in Fig. 6. The second and third rows show the frame-wise pedal target and system output. Values close to 1 indicate “on” states and values close to 0 indicate “off” states. The fourth and fifth rows show the regression onset targets and outputs. The sixth and seventh rows show the regression offset targets and outputs. Fig. 10 shows that our pedal transcription system performs well on the 5-second audio clip example.

### E. Error Analysis

Error analysis was carried on a few transcribed pieces from real recordings. We observe that most false positives come from octave errors. Usually, the harmonics of notes are recognized as false positive notes. Other false positive errors include false positives with short durations of less than 15 milliseconds and false positive of repeated notes. There are also a small number of false positives that do not have an obvious explanation. For false negative (missing) errors, we observe that sometimes the higher notes of octaves can be missed. In addition, there are a small number of bass notes ignored by our system. We also observe degraded performance when a piano is out of tune or the qualities of recording devices are low.

## V. CONCLUSION

We propose a high-resolution piano transcription system by regressing the precise onset and offset times of piano notes and pedals. At inference, we propose an analytical algorithm to

calculate the precise onset and offset times. We show that our proposed system achieves a state-of-the-art onset F1 score of 96.72% in piano note transcription, outperforming the onsets and frames system of 94.80%. We show that our system is robust to the misaligned onset and offset labels. In addition, we investigate evaluating piano transcription systems with different onset and offset tolerances that were not evaluated in previous works. As far as we know, we are the first to evaluate pedal transcription on the MAESTRO dataset and achieves a pedal event F1 score of 91.86%. One of the applications of our piano transcription system is the creation of the GiantMIDI-Piano dataset<sup>1</sup>. Other applications include piano performance analysis, genre analysis, etc. The limitation of our proposed systems includes the transcription results depend on the quality of audio recordings and the transcription system need to be modified for real-time applications. In the future, we will investigate multi-instrument transcription using our proposed high-resolution transcription system.

## VI. ACKNOWLEDGEMENT

We thank Mr. Mick Hamer for providing insightful analysis of our transcription system on several music pieces recorded on a Bosendorfer SE. We thank Prof. Gus Xia for providing a Yamaha Disklavier for playing back several transcribed MIDI files. We thank Mr. Hanying Feng for discussions on piano recording techniques.

## REFERENCES

- [1] C. Raphael, "Automatic transcription of piano music," in *International Society for Music Information Retrieval (ISMIR)*, 2002.
- [2] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 779–788.
- [3] E. Benetos, S. Dixon, Z. Duan, and S. Ewert, "Automatic music transcription: An overview," *IEEE Signal Processing Magazine*, vol. 36, no. 1, pp. 20–30, 2018.
- [4] E. Benetos, S. Dixon, D. Giannoulis, H. Kirchhoff, and A. Klapuri, "Automatic music transcription: challenges and future directions," *Journal of Intelligent Information Systems*, vol. 41, no. 3, pp. 407–434, 2013.
- [5] B. Li and Z. Duan, "An approach to score following for piano performances with the sustained effect," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 12, pp. 2425–2438, 2016.
- [6] B. Niedermayer and G. Widmer, "A multi-pass algorithm for accurate audio-to-score alignment," in *International Society for Music Information Retrieval (ISMIR)*, 2010, pp. 417–422.
- [7] Z. Duan and B. Pardo, "Soundprism: An online system for score-informed source separation of music audio," *IEEE Journal of Selected Topics in Signal Processing*, vol. 5, no. 6, pp. 1205–1215, 2011.
- [8] E. Scheirer, "Using musical knowledge to extract expressive performance information from audio recordings," in *IJCAI Workshop on Computational Auditory Scene Analysis*, 1995, pp. 153–160.
- [9] S. Dixon, "On the computer recognition of solo piano music," in *Proceedings of Australasian Computer Music Conference*, 2000, pp. 31–37.
- [10] M. Marolt, "Transcription of polyphonic piano music with neural networks," in *Mediterranean Electrotechnical Conference. Information Technology and Electrotechnology for the Mediterranean Countries*, vol. 2, 2000, pp. 512–515.
- [11] G. E. Poliner and D. P. W. Ellis, "A discriminative model for polyphonic piano transcription," *EURASIP Journal on Advances in Signal Processing*, 2006.
- [12] V. Emiya, R. Badeau, and B. David, "Multipitch estimation of piano sounds using a new probabilistic spectral smoothness principle," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 6, pp. 1643–1654, 2009.
- [13] J. P. Bello, L. Daudet, and M. Sandler, "Automatic piano transcription using frequency and time-domain information," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 6, pp. 2242–2251, 2006.
- [14] S. A. Abdallah and M. D. Plumbley, "Polyphonic music transcription by non-negative sparse coding of power spectra," in *International Conference on Music Information Retrieval (ISMIR)*, 2004, pp. 318–325.
- [15] K. O'Hanlon and M. D. Plumbley, "Polyphonic piano transcription using non-negative matrix factorisation with group sparsity," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2014, pp. 3112–3116.
- [16] T. Cheng, M. Mauch, E. Benetos, and S. Dixon, "An attack/decay model for piano transcription," in *International Society for Music Information Retrieval (ISMIR)*, 2016.
- [17] T. Berg-Kirkpatrick, J. Andreas, and D. Klein, "Unsupervised transcription of piano music," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2014, pp. 1538–1546.
- [18] M. Marolt, "A connectionist approach to automatic transcription of polyphonic piano music," *IEEE Transactions on Multimedia*, vol. 6, no. 3, pp. 439–449, 2004.
- [19] A. Cogliati, Z. Duan, and B. Wohlberg, "Piano music transcription with fast convolutional sparse coding," in *IEEE International Workshop on Machine Learning for Signal Processing (MLSP)*, 2015, pp. 1–6.
- [20] J. Nam, J. Ngiam, H. Lee, and M. Slaney, "A classification-based polyphonic piano transcription approach using learned feature representations," in *International Society for Music Information Retrieval (ISMIR)*, 2011, pp. 175–180.
- [21] S. Böck and M. Schedl, "Polyphonic piano note transcription with recurrent neural networks," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2012, pp. 121–124.
- [22] S. Sigtia, E. Benetos, and S. Dixon, "An end-to-end neural network for polyphonic piano music transcription," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 5, pp. 927–939, 2016.
- [23] R. Kelz, M. Dorfer, F. Korzeniowski, S. Böck, A. Arzt, and G. Widmer, "On the potential of simple framewise approaches to piano transcription," in *International Society for Music Information Retrieval (ISMIR)*, 2016.
- [24] R. Kelz, S. Böck, and C. Widnaer, "Multitask learning for polyphonic piano transcription, a case study," in *IEEE International Workshop on Multilayer Music Representation and Processing (MMRP)*, 2019, pp. 85–91.
- [25] C. Hawthorne, E. Elsen, J. Song, A. Roberts, I. Simon, C. Raffel, J. Engel, S. Oore, and D. Eck, "Onsets and frames: Dual-objective piano transcription," in *International Society for Music Information Retrieval Conference (ISMIR)*, 2018.
- [26] C. Hawthorne, A. Stasyuk, A. Roberts, I. Simon, C. A. Huang, S. Dieleman, E. Elsen, J. Engel, and D. Eck, "Enabling factorized piano music modeling and generation with the maestro dataset," *International Conference on Learning Representations (ICLR)*, 2018.
- [27] J. Schroeter, K. Sidorov, and D. Marshall, "Softloc: Robust temporal localization under label misalignment," *Open Review at International Conference on Learning Representations (ICLR)*, 2019.
- [28] A. Elowsson, "Polyphonic pitch tracking with deep layered learning," *The Journal of the Acoustical Society of America*, vol. 148, no. 1, pp. 446–468, 2020.
- [29] A. Elowsson, "Modeling music: studies of music transcription, music perception and music production," *KTH Royal Institute of Technology*, 2018.
- [30] A. Gkiokas and V. Katsouros, "Convolutional neural networks for real-time beat tracking: A dancing robot application," in *International Society for Music Information Retrieval (ISMIR)*, 2017, pp. 286–293.
- [31] W. Goebel, "Melody lead in piano performance: Expressive device or artifact?" *The Journal of the Acoustical Society of America*, vol. 110, no. 1, pp. 563–572, 2001.
- [32] Goebel, Werner, "The role of timing and intensity in the production and perception of melody in expressive piano performance," *PhD Thesis*, 2003.
- [33] R. Kelz, S. Böck, and G. Widmer, "Deep polyphonic ADSR piano note transcription," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 246–250.
- [34] R. M. Bittner, B. McFee, J. Salamon, P. Li, and J. P. Bello, "Deep salience representations for F0 estimation in polyphonic music," in *International Society for Music Information Retrieval (ISMIR)*, 2017, pp. 63–70.

<sup>1</sup><https://github.com/bytedance/GiantMIDI-Piano>

- [35] K. Ullrich, J. Schlüter, and T. Grill, “Boundary detection in music structure analysis using convolutional neural networks.” in *International Society for Music Information Retrieval (ISMIR)*, 2014, pp. 417–422.
- [36] R. B. Dannenberg, “The interpretation of midi velocity,” in *International Computer Music Conference (ICMC)*, 2006.
- [37] B. Liang, G. Fazekas, and M. Sandler, “Piano sustain-pedal detection using convolutional neural networks,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 241–245.
- [38] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” *arXiv preprint arXiv:1912.01703*, 2019.
- [39] Q. Kong, Y. Cao, T. Iqbal, Y. Wang, W. Wang, and M. D. Plumbley, “PANNs: Large-scale pretrained audio neural networks for audio pattern recognition,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 2880–2894, 2020.
- [40] J. W. Kim and J. P. Bello, “Adversarial learning for improved onsets and frames music transcription,” in *International Society for Music Information Retrieval (ISMIR)*, 2019, pp. 670–677.
- [41] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *International Conference on Machine Learning (ICML)*, 2015, pp. 448–456.
- [42] V. Nair and G. E. Hinton, “Rectified linear units improve restricted boltzmann machines,” in *Proceedings of International Conference on Machine Learning (ICML)*, 2010, pp. 807–814.
- [43] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [44] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *International Conference on Learning Representations (ICLR)*, 2015.
- [45] C. Raffel, B. McFee, E. J. Humphrey, J. Salamon, O. Nieto, D. Liang, D. P. Ellis, and C. C. Raffel, “mir\_eval: A transparent implementation of common mir metrics,” in *In Proceedings International Society for Music Information Retrieval Conference (ISMIR)*, 2014.