# Homework 3: Automatic Polyphonic Piano Transcription

20213013 Jiho Kang

November 22th, 2021

## [Experiments and Results]

The experiment was conducted under the same conditions (given settings) for all models (baseline, rnn, crnn, onf) except for iteration. The results are shown in Tables 1,2.

| | | Baseline | RNN | CRNN | ONF |
|---|---|---|---|---|---|
| **Frame** | **Frame Precision** | $0.814 \pm 0.049$ | $0.781 \pm 0.081$ | $0.681 \pm 0.059$ | $0.632 \pm 0.062$ |
| | **Frame Recall** | $0.559 \pm 0.094$ | $0.232 \pm 0.061$ | $0.483 \pm 0.081$ | $0.394 \pm 0.090$ |
| | **Frame F1** | $0.658 \pm 0.067$ | $0.354 \pm 0.074$ | $0.562 \pm 0.067$ | $0.482 \pm 0.079$ |
| | **Onset Precision** | $0.833 \pm 0.025$ | $0.724 \pm 0.042$ | $0.814 \pm 0.035$ | $0.813 \pm 0.036$ |
| | **Onset Recall** | $0.660 \pm 0.117$ | $0.314 \pm 0.140$ | $0.601 \pm 0.127$ | $0.597 \pm 0.128$ |
| | **Onset F1** | $0.732 \pm 0.080$ | $0.422 \pm 0.132$ | $0.686 \pm 0.094$ | $0.683 \pm 0.095$ |
| **Note** | **Precision** | $0.986 \pm 0.007$ | $0.928 \pm 0.030$ | $0.972 \pm 0.014$ | $0.967 \pm 0.015$ |
| | **Recall** | $0.720 \pm 0.117$ | $0.348 \pm 0.147$ | $0.659 \pm 0.129$ | $0.655 \pm 0.130$ |
| | **F1** | $0.827 \pm 0.079$ | $0.490 \pm 0.149$ | $0.779 \pm 0.096$ | $0.775 \pm 0.097$ |
| | **Overlap** | $0.575 \pm 0.054$ | $0.370 \pm 0.085$ | $0.492 \pm 0.063$ | $0.437 \pm 0.086$ |
| **Note with Offset** | **Precision** | $0.492 \pm 0.115$ | $0.275 \pm 0.148$ | $0.428 \pm 0.114$ | $0.379 \pm 0.135$ |
| | **Recall** | $0.363 \pm 0.119$ | $0.107 \pm 0.082$ | $0.291 \pm 0.102$ | $0.256 \pm 0.110$ |
| | **F1** | $0.415 \pm 0.117$ | $0.149 \pm 0.104$ | $0.344 \pm 0.106$ | $0.303 \pm 0.119$ |
| | **Overlap** | $0.863 \pm 0.058$ | $0.787 \pm 0.083$ | $0.850 \pm 0.074$ | $0.836 \pm 0.074$ |

**Table 1.** Precision, Recall, and F1 results on test dataset (iteration value is 10,000)

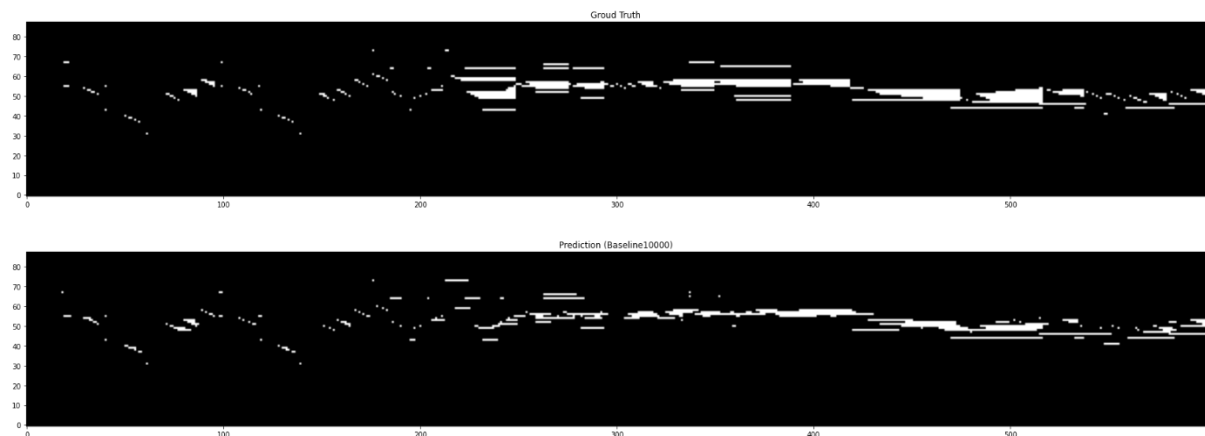| | | Baseline | RNN | CRNN | ONF |
|---|---|---|---|---|---|
| **Frame** | **Frame Precision** | $0.785 \pm 0.051$ | $0.740 \pm 0.064$ | $0.717 \pm 0.049$ | $0.657 \pm 0.059$ |
| | **Frame Recall** | $0.609 \pm 0.097$ | $0.423 \pm 0.088$ | $0.636 \pm 0.076$ | $0.499 \pm 0.079$ |
| | **Frame F1** | $0.680 \pm 0.062$ | $0.533 \pm 0.079$ | $0.672 \pm 0.054$ | $0.565 \pm 0.065$ |
| | **Onset Precision** | $0.832 \pm 0.027$ | $0.761 \pm 0.039$ | $0.823 \pm 0.035$ | $0.816 \pm 0.035$ |
| | **Onset Recall** | $0.679 \pm 0.114$ | $0.395 \pm 0.139$ | $0.686 \pm 0.110$ | $0.681 \pm 0.114$ |
| | **Onset F1** | $0.744 \pm 0.076$ | $0.508 \pm 0.121$ | $0.745 \pm 0.078$ | $0.739 \pm 0.080$ |
| **Note** | **Precision** | $0.988 \pm 0.006$ | $0.947 \pm 0.024$ | $0.974 \pm 0.013$ | $0.970 \pm 0.015$ |
| | **Recall** | $0.741 \pm 0.113$ | $0.440 \pm 0.147$ | $0.757 \pm 0.111$ | $0.751 \pm 0.114$ |
| | **F1** | $0.842 \pm 0.075$ | $0.587 \pm 0.134$ | $0.848 \pm 0.075$ | $0.842 \pm 0.078$ |
| | **Overlap** | $0.597 \pm 0.052$ | $0.459 \pm 0.050$ | $0.588 \pm 0.047$ | $0.495 \pm 0.062$ |
| **Note with Offset** | **Precision** | $0.514 \pm 0.115$ | $0.327 \pm 0.112$ | $0.525 \pm 0.091$ | $0.432 \pm 0.114$ |
| | **Recall** | $0.389 \pm 0.119$ | $0.157 \pm 0.088$ | $0.409 \pm 0.099$ | $0.335 \pm 0.109$ |
| | **F1** | $0.440 \pm 0.117$ | $0.208 \pm 0.099$ | $0.458 \pm 0.094$ | $0.376 \pm 0.110$ |
| | **Overlap** | $0.870 \pm 0.058$ | $0.823 \pm 0.079$ | $0.871 \pm 0.068$ | $0.852 \pm 0.072$ |

**Table 2.** Precision, Recall, and F1 results on test dataset (iteration value is 20,000)

# [Discussion]

1. Visualize at least one sample of your prediction (onset and frame) in the piano roll format



**Figure 1.** Visualization of onset (ground truth and prediction). A Part of the file 'MIDI-Unprocessed_SMF_17_R1_2004_03-06_ORIG_MID--AUDIO_20_R2_2004_12_Track12_wav' file'
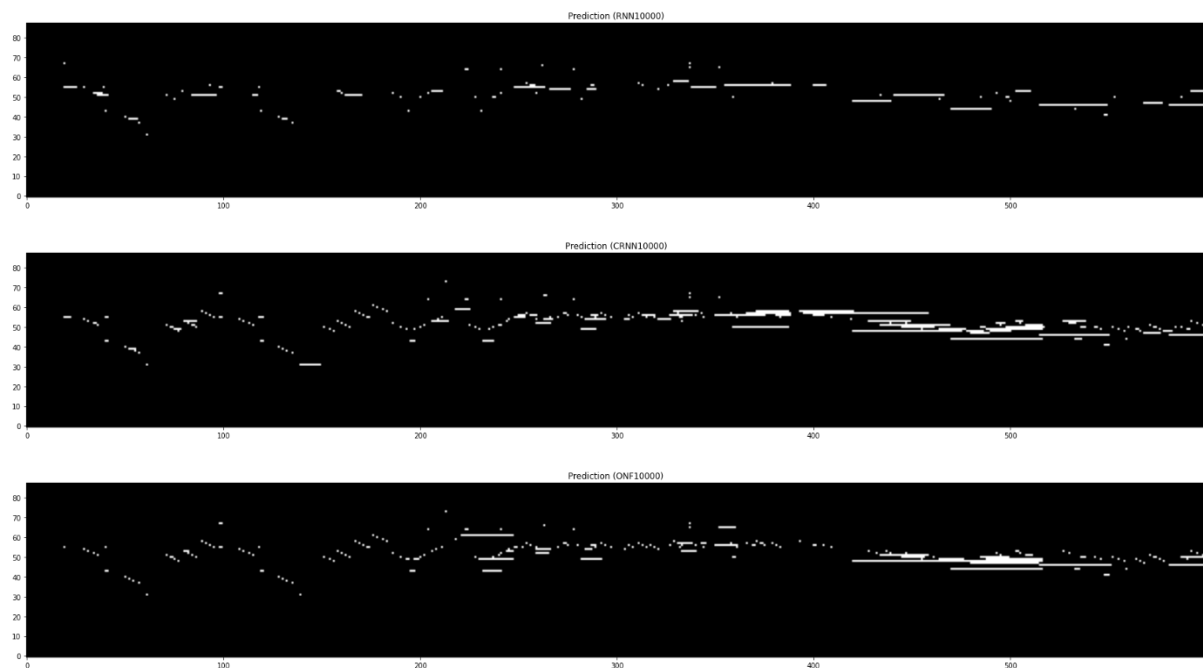
**Figure 2.** Visualization of frame (ground truth and prediction). A Part of the file 'MIDI-Unprocessed_SMF_17_R1_2004_03-06_ORIG_MID--AUDIO_20_R2_2004_12_Track12_wav'

2. What kinds of errors did you observe?

    A. Are the predicted onsets and frames consistent with each other?

**Figure 3.** Visualization of predicted onset and frame respectively, the result of overlapping onset and frame (Baseline)
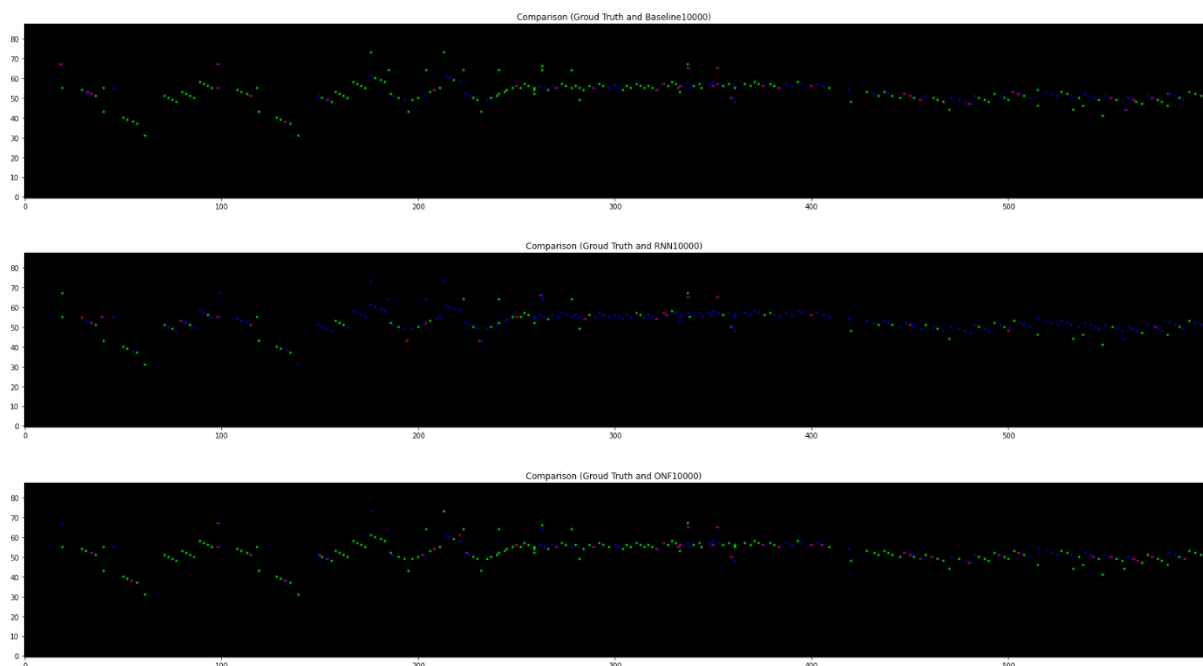


**Figure 4.** The result of overlapping onset and frame (RNN, CRNN, ONF)

As can be seen from Figures 3 and 4, all models showed that predicted onsets and frames consistent with each somewhat. One question is that the onsets appear again before the frame is turned off in addition to moment when the frame is turned on, as shown in figure 9 where the onset of the ground truth overlaps with the frame of the ground truth is displayed in purple. To distinguish these, the onset at the moment the frame is turned on is displayed in green, and the other onsets (can be interpreted as a re-onset) are displayed in red (figure 3,4,9).

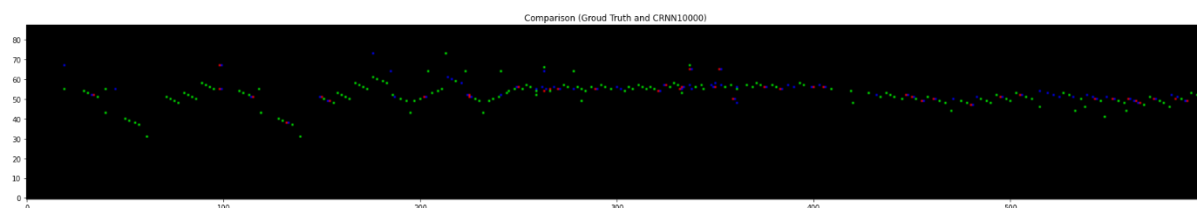B.  Compare them with the ground truth and analyze the errors in both frame-wise and note-wise perspective.

**Figure 5.** Visualization of onset (Comparison between ground truth and prediction). True Positive (Green), False Positive (Red), False Negative (Blue)
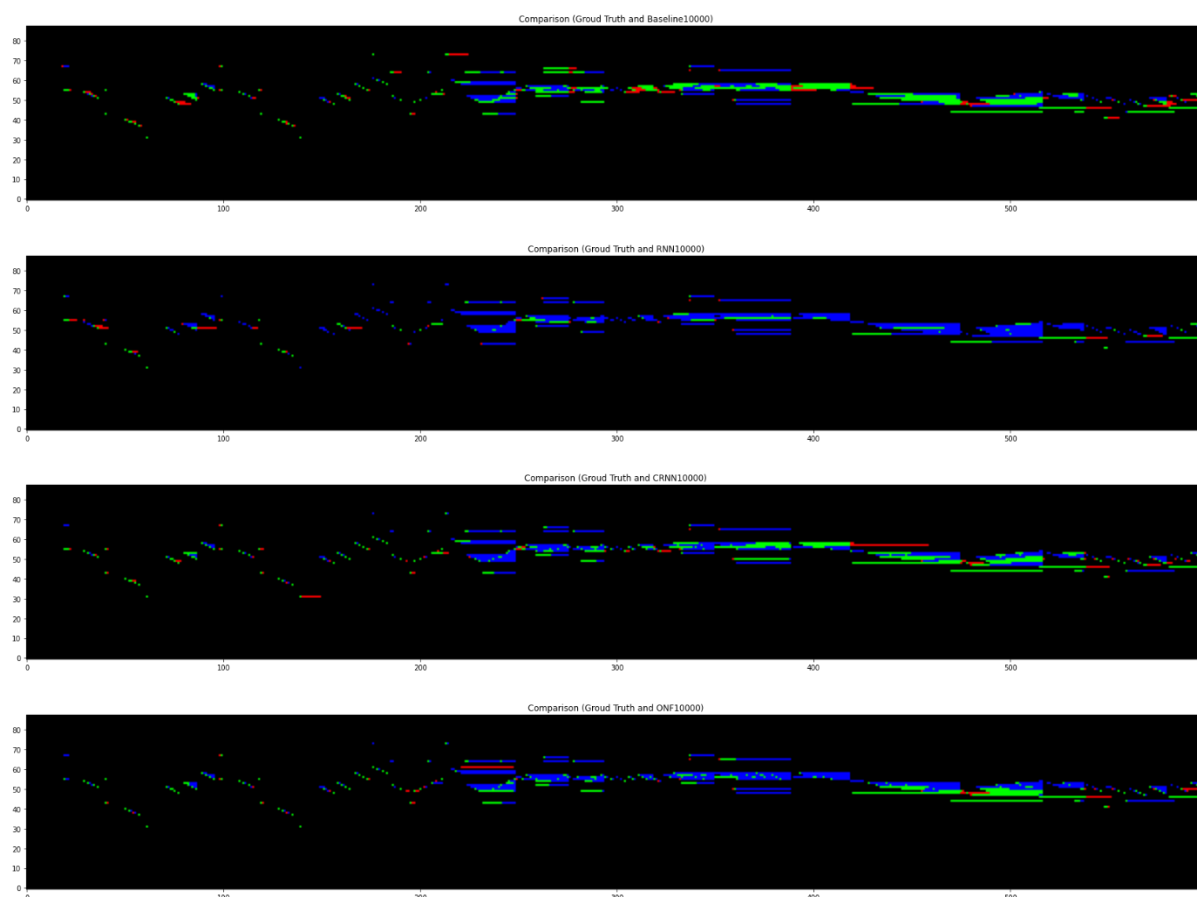


**Figure 6.** Visualization of frame (Comparison between ground truth and prediction). True Positive (Green), False Positive (Red), False Negative (Blue)

To compare predicted onsets and frames with the ground truth and analyze the errors in both frame-wise and note-wise perspective, the visualization was conducted in the way as figures 5 and 6. In figures, green pixels represents true positive, blue pixels represents false positive, and red pixels represents false. For both aspects' onset and frames, contrary to expectations, Baseline performed the best, CRNN and ONF performed similarly, and RNN performed the worst as shown in figure 5,6 and table 1 (iteration value is set to 10,000). One notable point is that the proportion of false negatives in frame prediction is high. In other words, models cannot predict multiple musical notes for each column on the horizontal axis. Therefore, it can be considered that the proposed model structures of this assignment do not capture the synchronicity of musical notes properly.

3. How would you improve the results?

The key point in this assignment is separation of onset state detection part in network. A

connection from the onset prediction to the input of frame prediction network restrict prediction framewise detector [1]. This is what empowers note state in a polyphonic piano music transcription task. According to [3], note-level accuracy is more important than the frame-level accuracy in terms of human musical perception. In line with this aspect, the following direction can be investigation about more note states such as onset, sustain, offset, re-onset, and even detection of the sustain pedal. To this end, we can make use of unified neural network architecture which can predict multiple note states using a soft-max output with a single loss function [2].

## Reference

[1] Curtis Hawthorne, Erich Elsen, Jialin Song, Adam Roberts, Ian Simon, Colin Raffel, Jesse Engel, Sageev Oore, and Douglas Eck. Onsets and frames: Dual-objective piano transcription. In Proceedings of the 19th International Society for Music Information Retrieval Conference, 2018.

[2] T. Kwon, D. Jeong, and J. Nam, "Polyphonic piano transcription using autoregressive multi-state note model," in Proceedings of the 21st International Society for Music Information Retrieval Conference (ISMIR), 2020.

[3] GCT634/AI613: Musical Applications of Machine Learning (Fall 2021) ppt slide, "Automatic Music Transcription"
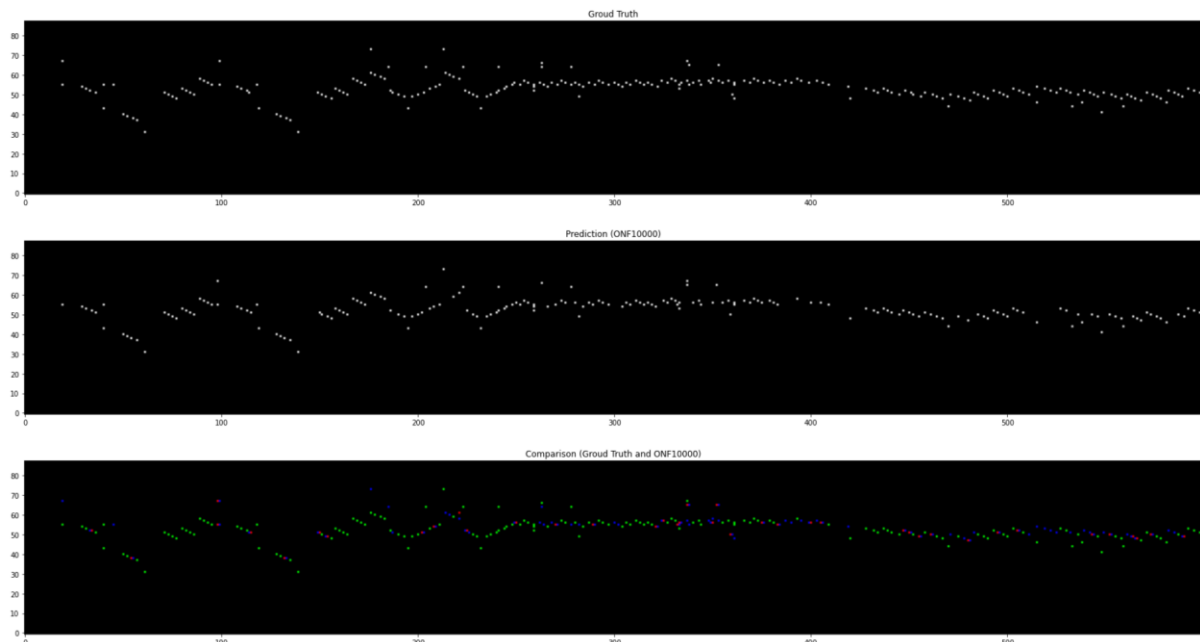
## Appendix



**Figure 7.** Visualization of onset (ground truth, prediction, and comparison)
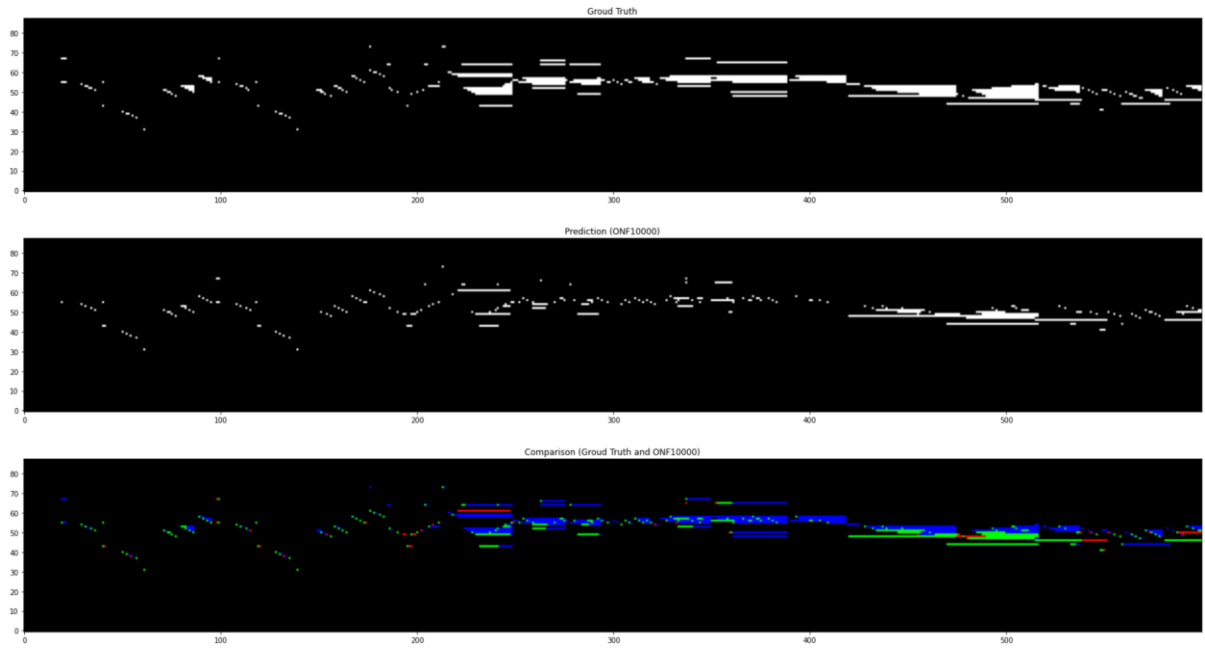
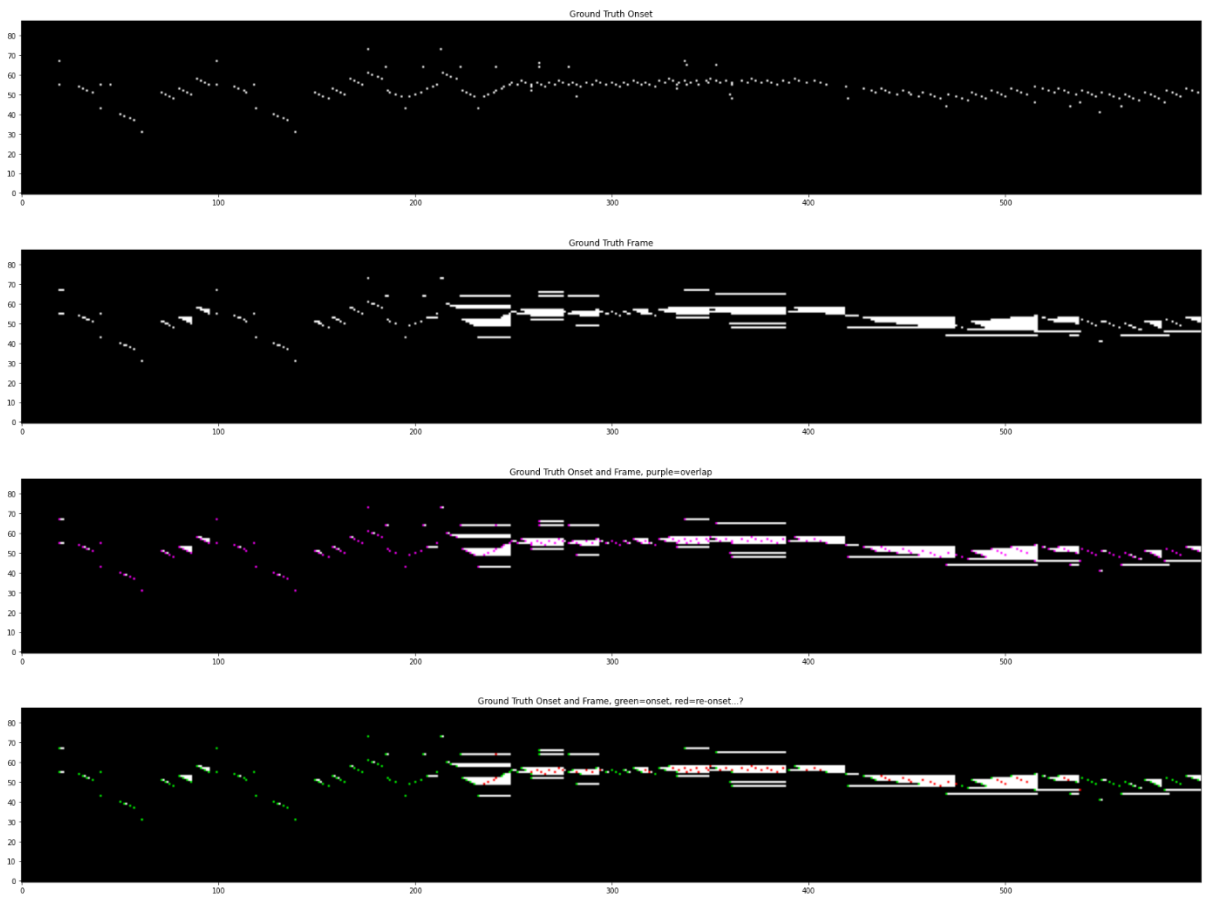**Figure 8.** Visualization of frame (ground truth, prediction, and comparison)



**Figure 9.** Visualization of ground truth onset and frame respectively, the result of overlapping onset and frame