

AUTOOSU: AUDIO-AWARE ACTION GENERATION FOR RHYTHM GAMES

Sihun Lee

Dasaem Jeong

Department of Art & Technology, Sogang University
Seoul, South Korea

ABSTRACT

Rhythm-based video games challenge players to match their actions with musical cues, turning songs into interactive experiences. The design of the game charts, which dictate the timing and placement of on-screen notes, are manually crafted by players and developers. With AutoOsu, we introduce a CRNN-based model for generating rhythm game charts for a given audio track, conditioned on an intended difficulty level. In previous studies, this task is often divided into two: onset detection, which determines timing points for notes; and action generation, where notes are distributed among a set of available keys. These two sub-tasks are typically handled with two separately trained models, and audio information is only given to the onset detection model. We instead jointly train the two recurrent layers who both receive audio information, which streamlines the training process and helps better utilize musical features.

1. INTRODUCTION

Rhythm games are a popular genre of modern video games. The gameplay of most rhythm games involves the player hitting specific keys at precise timings according to the notes that appear on the screen, the sequence of which is often called a *chart*. The charts are manually created by game developers or community members to follow the rhythmic and melodic structures of a song. By invoking a sense of moving along to the music, rhythm games provide players with a new and entertaining way of experiencing songs they like.

As a machine learning task, generating charts for rhythm games from a given audio input of music can be regarded as similar to music onset detection [1] and automatic music transcription [2]. However, a key difference is that there is no definitive answer for how one should chart for a given music track; a model needs to learn a wide range of idiosyncratic patterns in charts that are rooted in the physicality of how the games are played, the conventional note combinations used by the community, and how

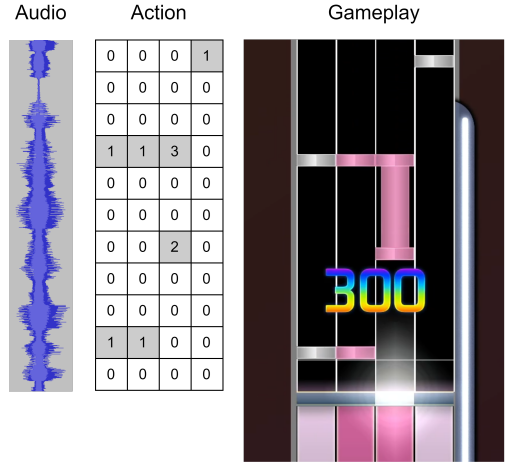


Figure 1. *osu! mania* gameplay. Notes fall down onto the line at the bottom, and the player has to hit the corresponding keys at the correct time.

chart creators tend to interpret different musical features of a song. In this sense, the task can also be regarded as conditional symbolic music generation.

In previous studies, this task is typically decomposed into two and handled with separate models. Donahue et al. [3] name the two sub-tasks *step placement* and *step selection*. In the former, the precise timing points of notes are determined—a process analogous to onset detection. In the latter, the notes are then distributed among a set of keys the player can hit. Liang et al. [4] adopt a similar 2-step approach composed of *timestep generation* and *action type generation*. For clarity, we define these two problems as *onset detection* and *action generation*.

In both studies mentioned above, the audio information is utilized by the onset detection model alone, and the action generation model is only conditioned on the time difference between the previous and the current note. In this research, we instead jointly train the two models and provide audio context to both modules, simplifying the training pipeline and fully utilizing the musical features extracted from audio tracks.

2. DATA

We focus on *osu! mania*, one of the game modes of the popular rhythm game "*osu!*", presented in Figure 1. Content for *osu!* is mainly produced by community members,



Number of songs	400 (16.4 hrs)
Avr. song length	148 secs
Number of charts	1,126
Notes / chart	676.54

Table 1. Dataset statistics

64 who create charts for songs and upload them to the game’s
65 database. Among the publicly available charts, we col-
66 lected 400 songs to compose the dataset, handpicking them
67 to maintain balance in genre and difficulty. Statistics of the
68 dataset are provided on Table 1.

69 3. METHOD

70 3.1 Feature Extraction

71 We extract raw audio tracks from charts in the dataset. To
72 preserve a wider range of low and high-level musical fea-
73 tures, we perform multiple timescale short-time Fourier
74 transforms in window lengths of 23ms, 46ms, and 93ms
75 [5]. We use a stride of 10ms, creating a grid of time frames
76 to which the inputs and outputs of all of the model’s com-
77 ponents are aligned.

78 Following [3], we compute rhythmic information for
79 each time frame: 1) *Beat number*, an integer that denotes
80 the beat in a measure that contains the time step; and 2)
81 *Beat phase*, representing the fraction of a beat at which the
82 time step occurs. For this, we assume that all audio tracks
83 are of consistent tempo and in a time signature of 4/4.

84 For action generation, we focus on the 4-key mode of
85 *osu! mania*. Each of the four keys can be assigned one
86 of the following actions at any time: *no note*, *normal note*,
87 *hold start*, and *hold end*. This results in a total of $4^4 = 256$
88 possible *action tokens* for each time step.

89 To condition chart generation on an intended difficulty
90 level, we collect the *star rating* of each chart, which is an
91 objective difficulty measure determined by the game’s in-
92 ternal logic. We limit the dataset to only contain charts of
93 difficulty levels lower than 4.0.

94 3.2 Model Architecture

95 As presented in Figure 2, the model comprises a stack
96 of convolution layers for processing audio, a bidirectional
97 Gated Recurrent Unit (GRU) [6] for onset detection, and an
98 auto-regressive unidirectional GRU for action generation,
99 which resembles an auto-regressive model for piano music
100 transcription [7]. The arriving spectrogram is forwarded
101 through the convolution stack with a gradual increase in
102 the channel dimension and then flattened along the channel
103 and frequency dimensions. While preserving the temporal
104 dimension, we concatenate the following tensors to the au-
105 dio representation: beat number embeddings, beat phase
106 embeddings, and difficulty projection, which is produced
107 by feeding a difficulty scalar into a multilayer perceptron.
108 The resulting concatenated tensor is used as input for both
109 GRUs.

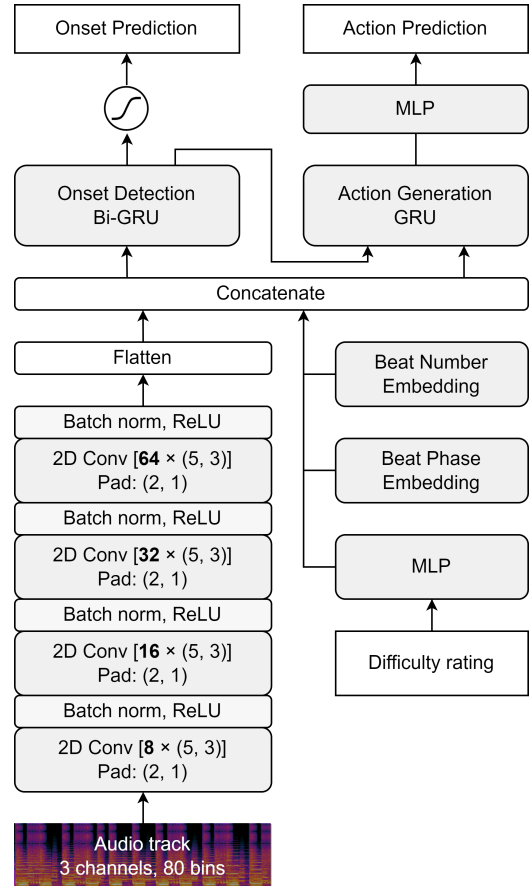


Figure 2. Overall model architecture

Additionally, the action generation GRU receives the
output of the onset detection GRU as part of its input. It
also receives the predicted action token from the previous
time step, making it the only autoregressive layer in the
model. Since the vast majority of ground-truth time steps
contain no notes, we utilize binary and multi-class focal
loss [8] for onset detection and action generation, respec-
tively, to mitigate class imbalance.

118 4. RESULTS

119 We compare the proposed model against a control model,
120 which only utilizes audio context during onset detection.
121 While the two models show no significant difference in
122 quantitative metrics, such as perplexity on the validation
123 set, we found that the generated charts by the proposed
124 model excel in aligning special patterns (hold notes, com-
125 pound notes) with salient musical events.

126 By inputting an audio track along with manual anno-
127 tation on tempo and the offset for the first downbeat, the
128 model can be used to perform inference with arbitrary
129 songs. For annotation, we use the tap-tempo feature in-
130 cluded in *osu!*, since it also calculates offset along with
131 tempo. Further examples and demos are provided in the
132 link¹. We also share the dataset, source code, and model
133 weights².

¹ <https://issyun.github.io/autoosu>

² <https://github.com/issyun/AutoOsu>

5. REFERENCES

- 135 [1] J. Schlüter and S. Böck, “Improved musical onset de-
136 tection with convolutional neural networks,” in *2014*
137 *IEEE International Conference on Acoustics, Speech and*
138 *Signal Processing (ICASSP)*. IEEE, 2014, pp. 6979–
139 6983.
- 140 [2] S. Sigtia, E. Benetos, and S. Dixon, “An end-to-end
141 neural network for polyphonic piano music transcrip-
142 tion,” *IEEE/ACM Transactions on Audio, Speech, and*
143 *Language Processing*, vol. 24, no. 5, pp. 927–939,
144 2016.
- 145 [3] C. Donahue, Z. C. Lipton, and J. McAuley, “Dance
146 dance convolution,” in *International conference on ma-*
147 *chine learning*. PMLR, 2017, pp. 1039–1048.
- 148 [4] Y. Liang, W. Li, and K. Ikeda, “Procedural content gen-
149 eration of rhythm games using deep learning methods,”
150 in *Entertainment Computing and Serious Games: First*
151 *IFIP TC 14 Joint International Conference, ICEC-*
152 *JCSG 2019, Arequipa, Peru, November 11–15, 2019,*
153 *Proceedings 1*. Springer, 2019, pp. 134–145.
- 154 [5] P. Hamel, Y. Bengio, and D. Eck, “Building musically-
155 relevant audio features through multiple timescale rep-
156 resentations,” 2012.
- 157 [6] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bah-
158 danau, F. Bougares, H. Schwenk, and Y. Bengio,
159 “Learning phrase representations using rnn encoder-
160 decoder for statistical machine translation,” in *Pro-*
161 *ceedings of the 2014 Conference on Empirical Meth-*
162 *ods in Natural Language Processing (EMNLP)*. As-
163 sociation for Computational Linguistics, 2014, p. 1724.
- 164 [7] T. Kwon, D. Jeong, and J. Nam, “Polyphonic pi-
165 ano transcription using autoregressive multi-state note
166 model,” in *Proc. of the 21th International Society*
167 *for Music Information Retrieval Conference (ISMIR)*,
168 2020.
- 169 [8] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár,
170 “Focal loss for dense object detection,” in *Proceedings*
171 *of the IEEE international conference on computer vi-*
172 *sion*, 2017, pp. 2980–2988.