# Toward a High Performance Piano Practice Support System for Beginners

Shota Asahi*, Satoshi Tamura*, Yuko Sugiyama† and Satoru Hayamizu*
*Gifu University, Gifu, Japan
E-mail: asahi@asr.info.gifu-u.ac.jp, tamura@info.gifu-u.ac.jp, hayamizu@gifu-u.ac.jp
†Chubu Gakuin College, Gifu, Japan
E-mail: ysugiyama@chubu-gu.ac.jp

*Abstract*— In piano learning, it is difficult especially for beginners to judge by themselves whether their musical performances are appropriate in terms of rhythm and melody. Therefore, we have been developing a piano practice support system, which enables piano beginners to conduct independent practice without their instructors. In this paper, we propose the system with the aid of a deep learning technique: Long Short-Term Memory (LSTM). Our system accepts raw piano sounds, extracting performance information. From these information, we evaluate performance. We evaluated the scheme using actual beginners' performances, and found the proposed system achieved better than previous conventional methods. This paper also presents an application employing our methods. Through subjective evaluation experiments for the proposed application, it turns out almost the all beginners found reflection points, and they maintained their motivation for independent practice.

## I    INTRODUCTION

In Japan, the ability of performing piano is one of the indispensable element for nursery teachers, and the lectures to acquire skills of piano performance are compulsory subject for students of nursery teacher training facility. On the other hand, among students there is a difference in piano performance experience before entering the facility. There is a limited amount of time for teaching students with less performing experience individually, and independent exercises are especially necessary for beginners of the piano performance to acquire sufficient performance skills. However, under the situations of independent practice in which any guidance by teachers cannot be got, it is difficult for the beginners to notice mistakes of the rhythm, melody or tempo when they are performing. This problem may have a negative impact on motivation and performance skill improvement. To solve these problems, researchers propose and develop a system to support independent practice. These researches are conducted from various perspectives, such as comparing the performance of students and teachers [1], visualizing practice situations [2]. Some methods applied for practice system [3].

We have researched to develop the piano practice support system which support beginners' independent practice from the perspective of music analysis. Our goal is to develop the web application software which can support independent practice by analyzing and evaluating their performance in real environments. We proposed the method of automatic evaluation that analyzes the piano performance recorded by the performer and evaluate the correctness of rhythm and melody [4]. We had used raw music signal data for performance data to make obtaining recording data easily from arbitrary practice environments. Music analysis from raw data with signal processing and machine learning have been researched by a considerable number of researchers, by applying the mechanism of speech recognition or other field's technique of signal processing to deal musical dimension such as melody, harmony, rhythm, etc. [5]. In our research, we had used the signal processing techniques such as Fast Fourier Transform (FFT) and power spectrum extraction.

To propose a high-performance system, we apply deep learning techniques for performance information extraction. Deep learning mechanisms, such as Recurrent Neural Networks (RNNs), which are usually used for time series data, have been used in the field of music analysis such as onset detection [6]. In our research, we use Long Short-Term Memory (LSTM) [7] to extract performance information such as each note's "pitch" and "timing" from the performance. In this paper, we evaluated our proposed information extraction method: whether our method has enough performance for accuracy evaluation. We also compared extraction performance with conventional methods to confirm whether it is appropriate to use an LSTM model for performance information extraction.

Finally, we implement ed proposed system into the practice support application and evaluated the effectiveness of system and application for the beginners' reflection and motivation.

## II    PIANO PRACTICE SUPPORT APPLICATION

One of our goals are to implement the performance information extraction method and performance evaluation method to a piano practice support application. Performance evaluation is done in the following order.

① Enter identification information of performer
② Capture recorded performance data (Fig.1)
③ Extract performance information
④ Evaluate the performance based on the extracted performance information
⑤ Visualize the evaluation results (Fig.2)

After inputting his / her identification information, the performer uploads the recorded performance to the system. After the upload is completed, the system extracts the timing
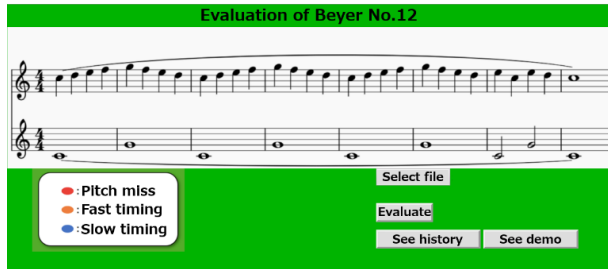
Fig.1. The screenshot of application's top page. Learners send recorded performance data from this page, and also can jump to the demo performance page and performance history page from this page.
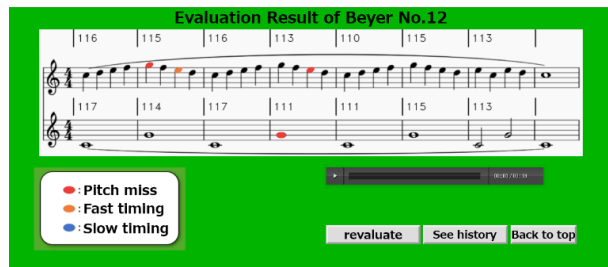


Fig.2. The screenshot of evaluation result page. Learners can feedback their performance with replaying recorded audio.

and pitch from the performance data and evaluates the melody, tempo and rhythm. After evaluation processing is completed, the evaluation result is plotted on the score image and displayed on the screen together with the recorded performance data. It is possible to perform performance feedback while listening to the actual performance. In the musical score image on which the performance evaluation result is plotted, a mistake in the melody / rhythm is presented by coloring the target note with red color when the musical interval is earlier than the correct one, yellow color when the performance rhythm is earlier than the correct one, and blue color when the performance rhythm is slower than the correct one. In the tempo of each measure, the numerical value of the tempo of the corresponding measure is set to "0" when the corresponding measure has replayed or not played note.

in this application we can see the history of performance evaluation. This function is not only for learners to see their improvement during independent practice, but also for teachers to check learners' progress and frequency of independent practice, and it can be used for reference of coaching.

### III    PIANO PRACTICE SUPPORT SYSTEM

The objective of this research is to automatically evaluate performance data recorded by web application implementing the system and show them and evaluation results to learners. We introduce the outline and the flow of the piano practice support system/application proposed for the purpose of helping beginners' independent practice.
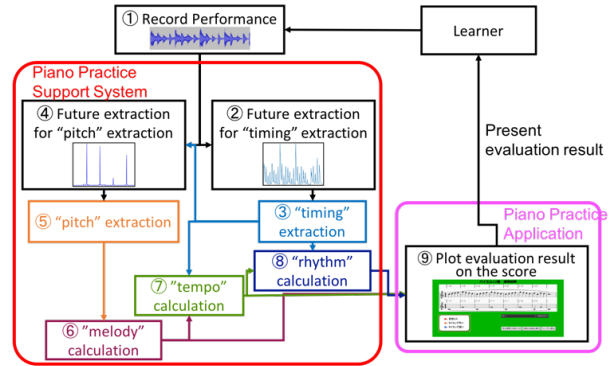


Fig.3. Flow of Piano practice support system.

#### A.    Performance Information and Evaluation Metrics

In order to satisfy this objective, we defined performance information in piano performance and metrics called evaluation metrics. In our research, we focus on the correctness of "pitch" flow and the correctness of "timing" interval between performance sounds. Two information, "timing" and "pitch" are used as the performance information extracted from the piano performance. As evaluation metrics, we use three types of metrics: "melody", "rhythm" and "tempo"; "melody" is a metric calculated from the flow of pitch of the performance sound from the beginning of the music to the end; "rhythm" is calculated from "timing", and is a metric for evaluating whether the sounds are performed with correct intervals; "tempo" is a metric that evaluates whether the song was being performed at the correct speed; "tempo" is calculated for each measure, used to check the speed change in the music performed. We designed the proposed system assuming that independent practices could be made at any location. For this reason, performance evaluation targets are not limited to musical score information recorded using MIDI but musical performance data in which raw sounds are recorded. Therefore, to calculate three kinds of evaluation metrics, it is necessary to acquire performance information such as the "pitch" and "timing".

#### B.    Flow of the System

The outline of the system and the transition to the application are shown in Fig. 3. The flow of the proposed system is as follows. First of all, we acquire performance data recorded by learners (①). Next, we extract performance information (②, ④). Then we predict and extract "timing" and "pitch" from a pre-learned LSTM prediction model ( ③, ⑤ ). After the extraction, we evaluate the performance with calculating evaluation metrics from extracted performance information. We evaluate melody (⑥), then tempo (⑦), and rhythm (⑧). After the extraction / evaluation, the system presents the results to the learner on the implemented piano learning application (⑨). As shown in Fig.2, the evaluation results of melody and rhythm are presented by giving colors to notes, and tempo results are displayed above each measure.

## C. Performance Information Extraction

Our system learns and extracts "timing" and "pitch" from the recorded music using an LSTM model. If more performance data can be used for learning, deep learning architecture can be extracted with a slight difference in the timing between the parts that appears during the performance, the influence of the harmonic overtone and the individual difference of the performance reduced [6]. For this reason, Bidirectional LSTM architecture is adopted as a prediction model in our system. Performance information extraction was performed separately in "timing" and "pitch". To design extraction models more dependent on each music, we construct a different performance information extraction model for each music and construct a pitch extraction model by the number of parts of the performance music. Pitch extraction is performed only using information before and after each note performed, so we perform the timing extraction process first and then perform the pitch extraction process with timing information.

The input of the model is the power spectrum extracted from music with both "timing" and "pitch". Upon implementation, we used a power spectrum extracted by framing with a frame length of 25(ms) and a frame width of 10(ms) as an input of timing extraction model and by framing with a frame length of 500(ms) and a frame width of 250(ms), with the timing of the played note as the center as an input of "pitch" extraction model.

The output of prediction model is like an output of sequence-to-sequence model. The "timing" extraction model outputs every 25(ms) which is the frame shift width, and a (part + 1) dimensional element indicating whether performance is started at that time as output. The "pitch" extraction model output on each corresponding note is 38-dimensional vector corresponding on pitch number.

## D. Evaluation Method of the Performance

To present the result of evaluation to learners, we proposed the method of calculating evaluation indices from performance information. Fig.4 shows the evaluation index on the score.

"pitch" is arranged in the order of the notes ("performed melody"). Then compare the "performed melody" with the "correct melody" labeled for each music, and record whether each note is played with accurate "pitch". At the same time, we check whether the replayed or not played sound is exist or not. If there are notes that have not been replayed or not played, this also affects the next rhythm evaluation and tempo evaluation. Finally, "melody" evaluation is done by classifying each note of the "performed melody" into 4 types of items: "correct answer sound", "pitch miss sound", "re-performed sound" and "not-performed sound".

As a tempo representation, we use Beats Per Minute (BPM). Tempo is calculated separately on each measure and on each part, so that you can check the tempo increase / decrease between measures, the transition of tempo with progress of the music, and the variation between the right hand and the left hand. The tempo is calculated using the extracted timing and melody evaluation. We use melody evaluation to exclude  from the calculation target measure which has replayed / not played
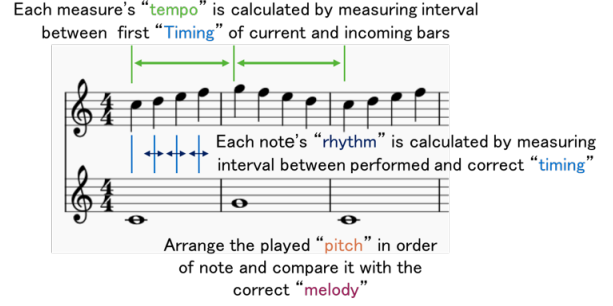


Fig.4. Relation between evaluation method and performance information extracted from performance data.

sound. The tempo $Tempo(n)$ of the n-th measure is calculated by equation (1), using "timing" of the first note of target measure $T_1^{(n)}$ and the next measure $T_1^{(n+1)}$. Since the power spectrum for "timing" extraction is conducted with a frame shift of 10ms, we convert the time between two notes into seconds for calculation

$$Tempo(n) = \frac{60[s]}{(T_1^{(n+1)} - T_1^{(n)}) \times 60/100} \quad (1)$$

To evaluate the rhythm, "correct timing" within each measure is derived from the tempo calculated first. For Beyer No.12, for example, 4 notes are in each measure. Therefore, "correct timing" at the $k^{th}$ ($k = 2,3,4$) sound is calculated is calculated by equation (1), using "timing" of the first note of target measure $T_1^{(n)}$ and "tempo" of belonging measure $Tempo(n)$. Can be expressed by the following equation (2). As same as "tempo" calculation, we don't perform rhythm evaluation if the n-th measure has replayed / not played sound.

$$CorrectTiming(n,k) = T_1^{(n)} + \frac{Tempo(n)}{4 \times 60}(k-1) \times 100 \quad (2)$$

After calculation the collect timing of each sound, we compare "correct timing" and extracted "timing" and check whether the "rhythm" of each sound is performed accurately, respectively. For judging accuracy and inaccuracy of rhythm, "1 beat $\times \pm 0.29(s)$" which was used as the evaluation value in on preliminary experiment. For example, in the case where the tempo of a bar in Beyer No.12 is BPM 120, since the length of one beat is 0.5(s), we assume the performed note as "correct rhythm" if "timing" is within $\pm 0.5 \times 0.29 = \pm 0.145(s)$. Finally, "rhythm" evaluation is done by classifying "performed sound" by 3 types of items: "correct answer", "fast rhythm", "slow rhythm", and "unevaluable sound" such as replayed / not played sound.

## IV  PERFORMANCE INFORMATION EXTRACTION EXPERIMENT

To evaluate our performance information extracting method, we conducted an experiment to evaluate the accuracy of performance information extraction using an extracted model trained from actual piano performance data.

Fig.5. Beyer No.12



Fig. 6. Beyer No.78

Table 1. The dimension of each layer and the type of hidden layer on proposed extraction architecture.

|  | Timing | Pitch |
|---|---|---|
| Input layer dim. | 40 | 600 |
| Hidden layer type | LSTM | BLSTM |
| Hidden layer dim. | 51 | 3 |
| Output layer dim. | 3 | 38 |

*A. Performance Data for Experiment*

The performance data used for the training of the model and the prediction for extracting the performance information was recorded using the actual initial scholar's piano performance. We recorded the performance in a small office and performed using YAMAHA upright piano. We used Beyer No.12 (Fig. 5) and the last section of No.78 (Fig. 6) for the experimental music. The performer is 18 junior college students, and all of performer are beginners. Each performer performed 2 music for 5 times respectively. The performance was recorded with a built-in microphone in the PC. Performance was cut out from the beginning of the music to the end by the recording software, and output as one performance data. Labeling of "timing" and "pitch" of each note was done manually. The recorded performance data of No. 12 was also used as learning data when implementing the application to be performed in the next chapter.

*B. Experimental Settings*

The purpose of the experiment is to investigate whether information extraction can be performed with sufficient accuracy for use as piano practice support system. To fulfill this purpose, we conducted two experiments. The first is to validate whether proposed system performs sufficient accuracy for application by learning and extracting performance information from actual recorded data. The second is to compare the accuracy of the extraction method to investigate the effectiveness of using LSTM for information extraction.

Table 2. Precision, Recall F-score of the test sets on information extraction using LSTM model on Beyer No.12.

| No.12 |  | Right | Left | Total |
|---|---|---|---|---|
| Timing | Precision | 0.9361 | 1.0000 | 0.9513 |
|  | Recall | 0.9982 | 1.0000 | 0.9986 |
|  | **F-score** | **0.9961** | **1.0000** | **0.9744** |
| Pitch | Precision | 1.0000 | 1.0000 | 1.0000 |
|  | Recall | 0.9983 | 1.0000 | 1.0000 |
|  | **F-score** | **0.9991** | **1.0000** | **0.9993** |

The proposed extraction model was implemented on Keras [9] with TensorFlow [8] backend. For evaluating performance, we used Precision, Recall, and F-score, which are generally used evaluation criteria [10].

*C. Performance Information Extraction*

Firstly, we performed experiments to extract the performance information from the actual performance using the proposed extraction method and investigate the extraction accuracy. The number of dimensions of the input layer, the hidden layer (1 layer) and the output layer of the performance information extraction model are shown in Table 1. We separate recorded data to train sets and test sets. From 18 performers' data, 14 performers' data used as model learning sets and the 4 performers' used as test sets. For "timing" extraction, we decide that if the difference between extraction and label is within ± 5 frames i.e. within ± 50ms, model extracted correct "timing". For "pitch", we decided that if same "pitch" extracted as label, model extracted correctly. Timing and pitch extraction accuracy is shown in Table 2. The results of 3 kinds of right hand, left hand, and both hand (as Total) are shown, respectively. The value of Precision, Recall and F-score became 0.95 or more in all conditions and confirmed that the performance information can be extracted with enough accuracy to be used as the performance evaluation of the piano practice support system. These results showed that the prediction model had sufficient extraction accuracy for incorporation into the application.

We also extracted the performance information for Beyer No.78. Beyer No.78 consists of 2 parts for right hand and 1 part for left hand. With our method, we extracted "timing" and "pitch" information of the upper part of right hand and the part of left hand accurately. On the other hand, the "timing" extraction accuracy of the lower part of right hand declined compare with other parts. This problem might come from the less number of scores performed in the lower part of right hand than the other part, so we'll consider about how to divide the music into parts.

*D. Comparison of Extraction Accuracy*

Second, we investigate whether it is appropriate to use LSTM model to extract performance information by the model without deep learning architecture. For the comparison, we use 2 method: threshold processing with both timing and pitch extraction [4], and the extraction method by Support Vector Machine (SVM) [11]. We used same features used in LSTM model for SVM model. We used F-score for comparison metric.

Table 3. F-score of information extraction of the test sets on information extraction under different models, LSTM, threshold processing model and SVM on Beyer No. 12.

| No.12 | | Right | Left | Total |
|---|---|---|---|---|
| Timing | LSTM | 0.9611 | 1.0000 | 0.9744 |
| | threshold | 0.2481 | 0.0909 | 0.2160 |
| | SVM | 0.3442 | 0.4838 | 0.3666 |
| Pitch | LSTM | 0.9991 | 1.0000 | 0.9993 |
| | threshold | 0.9508 | 0.2098 | 0.6770 |
| | SVM | 0.8707 | 0.4781 | 0.7162 |

Timing and pitch extraction accuracy comparison results are shown in Table 3. The annotations in the graph are the same as those in the previous section. The extraction timing extraction and the SVM compared with the threshold processing and the SVM. This is thought to be due to the fact that the model by LSTM can reduce the environmental noise, the difference in performance timing between notes and the influence of harmonics when timing extraction. In the case of actual recording, Unlike MIDI data, even when playing in a silent environment, ambient noise such as a keyboard depression sound or resonance sound entering slightly, even though it can separate each part by overtone removal etc. With this factor, SVM extractor can possible to extract the timing of the left hand with high accuracy, but the accuracy may decrease at the time of extracting the timing of the right hand that interferes with harmonics. In terms of being able to extract without being influenced by these factors, it is effective to use the model by LSTM in the timing extraction method is there. For pitch estimation, the F-score of the pitch extraction of the left-hand decreases in the thresholding process and the extraction method by the SVM. This is because the band of the scale of the bass is higher than the treble. It is thought that the waveforms of the silent section and the sound section almost no change, in particular, so it is considered that factor can be extracted more accurately without being influenced by a slight difference due to these factors as well as the timing. It is effective to use the model by LSTM in the pitch extraction method.

## V    APPLICATION VALIDATION EXPERIMENT

In this chapter, we describe the experiment validating the effectiveness of the proposed system and application protocol by having the beginner use the proposed application.

### A.    Experiment for Application Evaluation

To validate the effect of the proposed application on the learning efficiency of the beginner performer and the motivation to learning, we conducted an evaluation experiment with 10 actual beginner performers as subjects. The subject experiment was divided into 3 tests per person. The procedure of the experiment is as follows.

①    At the beginning of the experiment, explain the flow of the experiment and the score and fingering of Beyer No. 12.
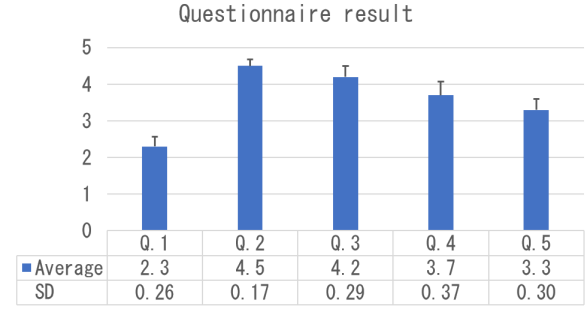②    Subjects practice voluntarily Beyer No.12 for 10 minutes.



Fig.7. The result of application evaluation experiment.

③    After the independent practice, the subject record the performance, browsing the evaluation by the application, and checking points of reflection and improvement points.
④    After three tests, subjects answer the questionnaire.

Questionnaire question items are shown in Table 4. Question 1 to question 5 are evaluated by the five-point scale [12], question 6 is an answering method by free description. The evaluation and consideration of the application was carried out by investigating the transition of the evaluation result in addition to the response of the questionnaire.

### B.    Result and Discussion

Among the questionnaires, the results of the answers of question 1 to question 5 are shown in Fig.7. Also, among the performance evaluation results of 10 subjects, the transition o performance evaluation result of subject No. 6 who answered "Easy" to question 1 is shown in Fig.8, the subject who answered "Cannot say either" to question 1 is shown in Fig.9. From the performance evaluation result, subject No. 6 didn't made mistakes in all evaluation. On the other hand, subject No. 9 made a mistake up to the second time, and in the third it became possible to play accurately. This result showed that the degree of difficulty exercised by the test subject himself felt correlated with practice progress. Proposed application is effective in that it can visually check the status of learning progress for each learner by using it for independent practice and also can be expected to be utilized for efficient teaching. From the answers from question 2 to question 6, consider interest and motivation for piano learning. For question 2, all the subjects answered "more than a little" or more and can feedback their own performance while listening to the recording. On the other hand, there are variations in the evaluation results for question 3 over question 2. In question 6, in the free description column, "Since I was able to feel my improvement by visualization", "Because I got a mistake in playing and I was able to practice intensely", I also found out the point of reflection and improvement It has been shown that it leads to the maintenance of motivation. However, there were also responses such as "There was a gap in my experience and evaluation results", "I did not have the ability to improve", I could not feel improvement and affect the reduction of motivation. In addition, all the subjects in this experiment were

Table 4.  Questionnaire items

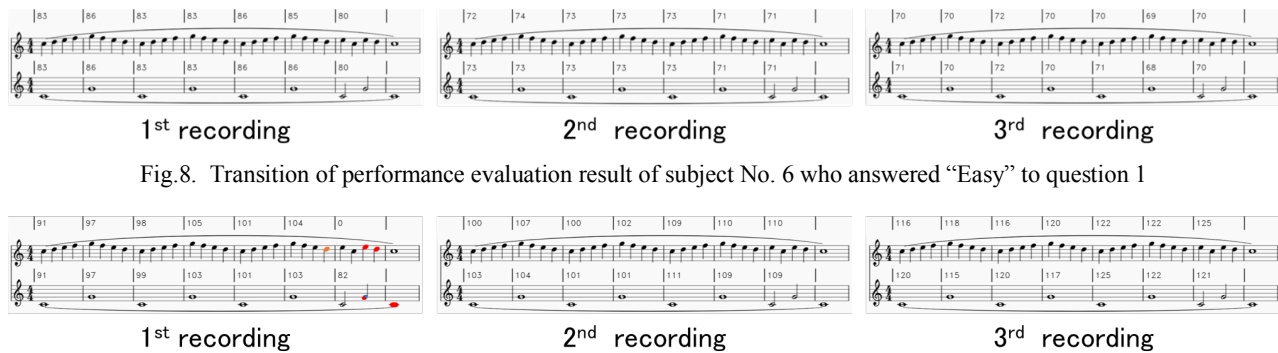| |
|---|
| Q.1 How was performance difficulty of Beyer.12? <br> 　1. Easy,  2. Not so difficult,  3. Cannot say either,  4. Little difficult,  5. Difficult |
| Q.2 Did you find the reflection point on your performance by performance visualization? <br> 　1. No,  2. Only a little,  3. Cannot say either,  4. Slightly yes,  5. Yes |
| Q.3 Did you feel that your performance improved with performance visualization? <br> 　1. No,  2. Only a little,  3. Cannot say either,  4. Slightly yes,  5. Yes |
| Q.4 Did you get interested in piano performance by using this application? <br> 　1. No,  2. Only a little,  3. Cannot say either,  4. Slightly yes,  5. Yes |
| Q.5 Did you think you'll continue to practice piano with this application? <br> 　1. No,  2. Only a little ,  3. Cannot say either,  4. Slightly yes,  5. Yes |
| Q.6 Free description |



Fig.8.  Transition of performance evaluation result of subject No. 6 who answered "Easy" to question 1



Fig.9.  Transition of performance evaluation result of subject No. 9 who answered "cannot say either" to question 1

students who don't have experience of piano performance, but six of them answered "slightly spoken" for question 4, but for question 5, It is understood that the answer value is decreasing. From these answers, it was shown that the use of this application may lead the first scholar to interest in piano learning and lead to motivation maintenance. Meanwhile, there is also a problem that improvement cannot be realized using the application, and it is necessary to make improvements such as presenting advice to improve learners and simplifying evaluation results.

## VI  CONCLUSION

We constructed a piano practice support system that enables beginner performers to feedback their own performance with themselves by analyzing and evaluating recorded raw sound. We conducted the experiment to evaluate the effect of our performance information extraction method using LSTM model, in most cases the F-score becomes 0.95 or more, and for Beyer 12 it is sufficient extraction accuracy to build the application It has been shown. We also conducted the experiment to evaluate the effect on piano practice with a web application implementing our proposed system. From the result, we found that all subjects understood the reflection points of the performance by feedback, and the application have some effect for improvement of performance and maintaining

motivation. Moreover, by using the application, it can be expected to improve efficiency of teaching piano by teachers.

For the future task, to improve the evaluation accuracy, we'll modify the performance information extraction method with two points, feature extracted from performance and part. About feature, we test other feature extraction method such as "contextual spectrums" [13]. About part, we found from experiment that when the number of part increase and each part has different number of scores, the extraction accuracy decreased on the part which has least number of scores. We'll test other part division method and compare extraction accuracy. For the application, we will conduct experiments such as comparing the case of learning using the application and the case of learning without using it and conducting experiments such as having the practiced performer use it, to discuss furthermore about the effect of the application. Moreover, we'll also add functions for maintaining practice motivation, such as improving the contents of the result display.

REFERENCES

[1]  K. Tanaka, T. Suzuki, M. Tsukamoto, "The Tendency of the MIDI data of a Piano Student and a Teacher ―From a Piano Teacher's Viewpoint―, "IPSJ SIG Technical Report Music Information Science (MUS), 2014, pp.1-6.

[2]  K. Ueda, Y. Tanaka, K. Hirata, "Evaluation of a Piano Learning Support System Focusing on Visualization of Keying

Information and Annotation," IPSJ SIG Technical Report Music Information Science (MUS), 2015, pp.1-8.

[3] K. Ueda, Y. Tanaka, K. Hirata, "Design and Implementation of a Piano Learning Support System Focusing on Visualization of Keying Information and Annotation," IPSJ Journal, Vol.57, No.12, 2016, pp.2617-2625.

[4] S. Asahi, S. Tamura, S. Hayamizu, Y. Sugiyama, "Estimation of tempo, timing, and melody for piano practice support systems", The Journal of the Acoustical Society of America, 140.4, 2016, pp3429.

[5] M. Müller, D. P. Ellis, A. Klapuri, G. Richard, "Signal processing for music analysis," IEEE Journal of Selected Topics in Signal Processing, 5(6), 2011, pp.1088-1110.

[6] F. Eyben, S. Böck, B. Schuller, A. Graves, "Universal onset detection with bidirectional long-short term memory neural networks," Proceedings of International Society of Music Information Retrieval Conference (ISMIR), 2010, pp. 589-594.

[7] F. A. Gers, J. Schmidhuber, F. Cummins, "Learning to forget: Continual prediction with LSTM", Neural computation 12(10), 2000, pp.2451-2471.

[8] M. Abadi, et al, "TensorFlow: A System for Large-Scale Machine Learning", Proceedings of the 12th USENIX conference on Operating Systems Design and Implementation (OSDI), 2016, pp. 265-283.

[9] Keras Documentation, Retrieved May 23, 2018, from https://keras.io/ .

[10] Powers, M. W. David, "Evaluation: From Precision, Recall and F-Measure to ROC, Informedness, Markedness & Correlation", Journal of Machine Learning Technologies, 2 (1), 2011, pp.37-63.

[11] C. Cortes, V. Vapnik, "Support-vector networks. Machine learning", 20(3), 1995, pp.273-297.

[12] R. Likert, "A technique for the measurement of attitudes", Archives of Psychology of New York 22(140), 1932, pp.5-55.

[13] L. Jiao, et al, "Contextual Spectrum Inpainting with Feature Learning and Context Losses", 7th International Conference on Information Society and Technology, 2017, pp.293-298.