# Digitalizing educational workbooks and collecting handwritten answers for automatic scoring

Tomo Asakura [1], Hung T. Nguyen [2], Nghia T. Truong [2], Nam T. Ly [2], Cuong T. Nguyen [2], Hiroshi Miyazawa [1], Yoichi Tsuchida [1], Takahiro Yamamoto [1], Masamitsu Ito [1], Toshihiko Horie [1], Fumiko Yasuno [3], Tsunenori Ishioka [4], Kokoro Kobayashi [5], Ikuko Shimizu [2] and Masaki Nakagawa [2]

[1] *Wacom Co., Ltd., 2-510-1 Toyonodai, Kazo-shi, Saitama, 349-1148, Japan*
[2] *Tokyo University of Agriculture and Technology, 2-24-16 Naka-cho, Koganei-shi, Tokyo, 184-8588, Japan*
[3] *National Institute for Educational Policy Research, 3-2-2 Kasumigaseki, Chiyoda-ku, 100-8951, Japan*
[4] *The National Center for University Entrance Examinations, 2-19-23 Komaba, Meguro-ku, 153-8501, Japan*
[5] *Ohshima College, National Institute of Technology, Ohshima, Yamaguchi, 742-2193, Japan*

**Abstract**
This workshop paper presents the digitalization of a set of educational workbooks on electronic paper, collection of digital ink answers, and description of correct answer examples as a first step toward realization of automatic scoring of handwritten answers. InkML was used to record handwritten answers. The questions and correct answer examples are associated with the handwritten answers. Digitalizing the entire workbooks involves the challenges associated with processing various types of questions in addition to ordinary questions asking for textual and mathematical answers. This paper also reports the prototypes of handwriting recognizers and automatic scorers for Japanese, mathematics, and English answers as well as the prototypes of user interfaces for answerers and human scorers. The goal is to achieve automatic scoring of handwritten answers in educational workbooks and computer-based testing, as well as to provide input by handwriting for learning management systems.

**Keywords**
Handwritten answers, handwriting recognition, automatic scoring, electronic paper, computer-based testing.

## 1. Introduction

Because of the widespread use of personal computers and tablets, education has become increasingly digitalized recently. Intelligent textbooks provide learners with multimodal contents, navigation, personalization, and so on, which paper textbooks cannot provide [1]. However, questions and answers seem to be limited in digitalization.

Multiple-choice questions are often used rather than descriptive questions because the input method is a keyboard or a touch panel, and such questions can be scored unambiguously. On the other hand, descriptive questions asked in paper-based textbooks, workbooks, and examinations can show the

understanding and problem-solving abilities of the answerers. The cognitive load is reduced for answerers as well as questioners designing the questions. Scoring descriptive answers, however, is labor-intensive and time-consuming. Feedback to answerers is often delayed, which decreases the motivation to review the questions and the effects of the review. A concern also exists that multiple-choice questions may provoke guessing of answers rather than problem-solving.

With recent advances in handwriting recognition, automatic recognition and scoring of handwritten Japanese constructed response answers composed of 80–120 characters has become possible with almost the same accuracy as that of humans [2, 3]. Automatic scoring of handwritten mathematical answers is also being studied to reduce the burden on the scorers [4]. Therefore, the extent to which any type of conventional questions can be recognized and scored automatically should be verified. The first step is to digitalize educational workbooks and collect handwritten answers using real workbooks with the aim of developing automatic scoring. Once the technology is realized, conventional workbook exams can be digitalized, and the scope of the questions asked under learning management systems can be enhanced.

Electronic paper with an electronic pen can be utilized to display questions and to capture handwritten answers in the form of online trajectory patterns, i.e., digital ink. Digital ink records the process through which each answerer answers each question, which can be analyzed to determine how the answerer has arrived at the answer. This approach is useful if the answer is wrong. Digital ink can be converted into an image so that both digital ink recognition (online recognition) and image recognition (offline recognition) can be used to produce reliable recognition.

Databases of handwritten character patterns in many languages have contributed to the progress of handwritten character recognition. Thus, a database of handwritten answers in the form of digital ink for actual workbook questions by many answerers should be prepared to develop automatic scoring, which includes correct and various incorrect answers.

The remainder of this workshop paper is organized as follows. Section 2 describes the related work. Section 3 introduces the device and data format. Section 4 presents prototypes of automatic recognizers and scorers. Section 5 describes prototypes of user interfaces for answerers and human scorers. Section 6 draws the conclusions.

## 2. Related work

In the 2000s, the utilization of tablets in the educational field started with the development of tablets and similar devices [5, 6]. Among them, our initial prototype was reported to score handwritten English vocabulary tests automatically [7, 8]. The benefit of capturing handwritten answers in digital ink was demonstrated [9, 10]. Although the input is image, a combination of handwriting recognition and natural language parsing was reported to score handwritten essays automatically [11].

Recently, many studies have been conducted to develop intelligent learning platforms for students such as ASSISTments [12] and Cognitive Tutor [13]. They focused on math problems and required photos of the answers provided by the students so that the human or computer-based scorers could give a score for each answer. In particular, the ASSISTments team organized the MathNet competition, where the collected dataset consisting of digitized student work was used to develop the handwriting recognition and automatic scoring systems [14]. On the other hand, this paper presents an attempt to digitalize multiple types of questions for elementary grades in Japanese, mathematics, and English subjects.

This paper presents three main contributions: a process for collecting handwritten answers using electronic paper, a prototype system of handwritten answer recognizers and scorers and a prototype of user interfaces for answerer and scorers.

## 3. Collected data

With the goal of digitalizing the workbooks as described in Section 2, we prepared an electronic paper device. This device displays each page of the workbooks and captures handwritten answers in digital ink. The details are provided in Subsections 3.1–3.4.

### 3.1. Device for collecting handwritten answers

The electronic paper employed in this study is equipped with an electronic pen, an eraser, and a compass, as shown in **Figure 1**, with a high sampling rate of 480 Hz. It records not only the pen/eraser/compass-tip coordinates, but also the writing pressure, tilt, and how far the tip is from the surface (up to 10 mm). Its dimensions are 209 mm×157 mm, which are close to those of an A5 sheet, and its thickness is 3 mm

The electronic paper connects to a host device (PC, smartphone, etc.) to display a page in one of the workbooks with a paper-like appearance and enables the user to answer with the pen or the compass and erase the trajectory with the eraser, which is recorded in digital ink. The host device receives the recorded digital ink.
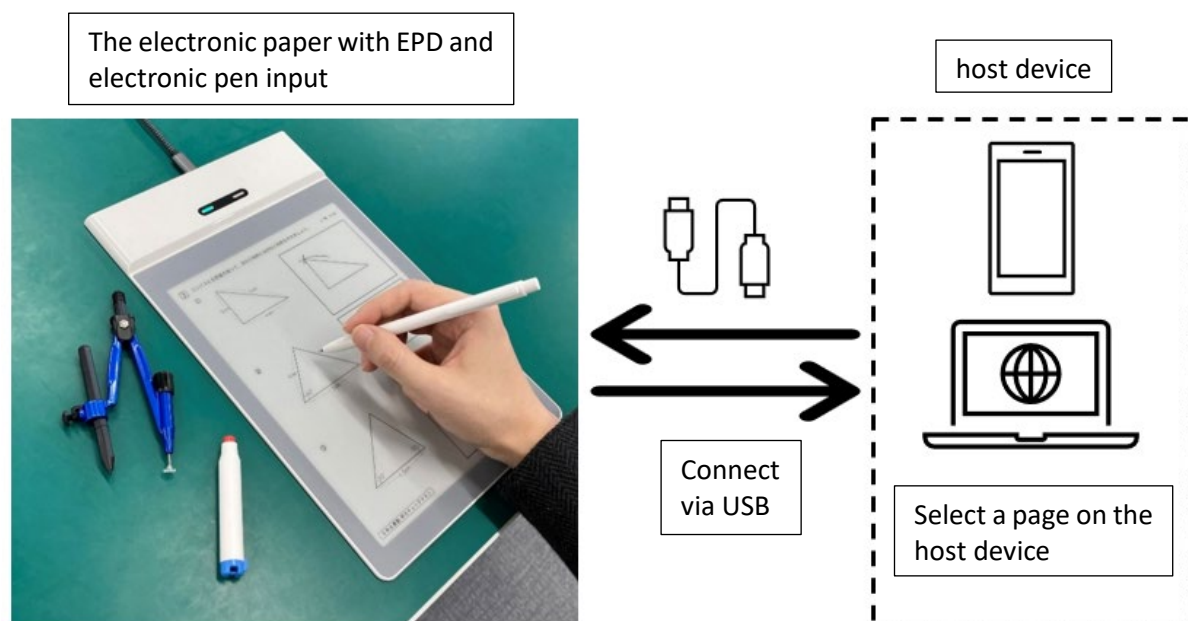


**Figure 1**: Electronic paper and connection to a host device.

### 3.2. Data format

For each page of the workbooks and for each answerer, the image, coordinates of the answer areas, correct answer examples, information about the answerer, and manual scoring results are recorded in XML. For the convenience of machine learning, we integrate the contents displayed on one screen (equivalent to one sheet of paper) into one file.

To annotate the question information, we use an <input> element with common attributes of inputId, the x- and y-coordinates of the top-left point of the input area, the width and height of the input area, the recognitionType for declaring the required recognizer for the input area, and the data for describing the details of the input area. To denote an expected answer, we use an <answer> element with two main attributes of inputId to match with inputId of the <input> element and expectedValue to present examples of the correct answer. In the future, we may include other attributes such as the condition to describe the rubric.

Sequences of sampled pen-tip/eraser-tip coordinates (digital ink) are recorded in the InkML format [15]. The coordinates of the pen/eraser/compass are recorded on the entire screen rather than in each answer area. The eraser does not erase the coordinates written by the pen but overwrites them in white ink of a specific thickness. All the coordinates from the pen and the eraser are recorded even outside the answer areas.

**Figure 2** shows the overall structure of the XML and InkML descriptions in a handwritten answers log file, and **Figure 3** provides an example of digital ink recorded in InkML and a correct answer description in XML.
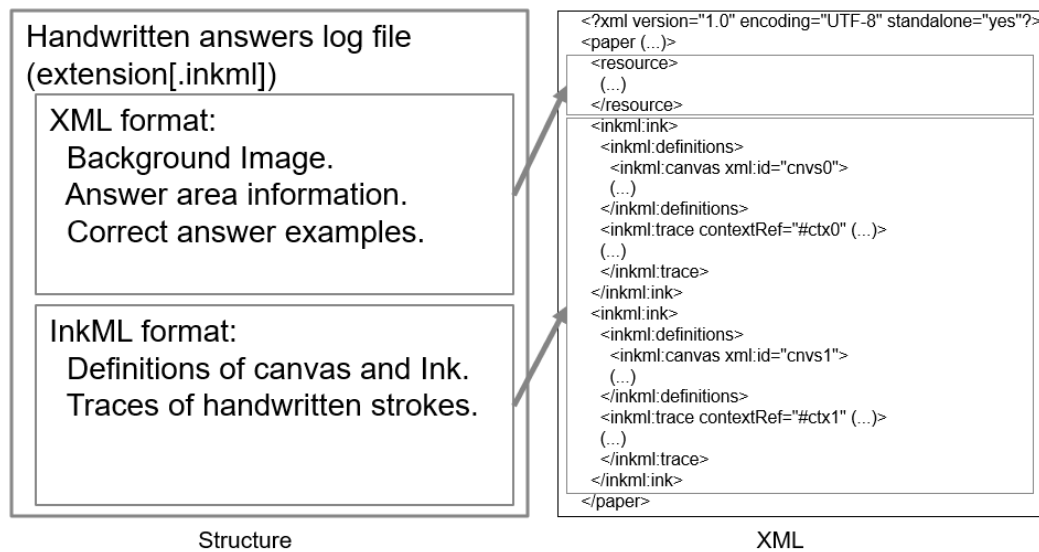


**Figure 2:** Data structure for handwritten answers.



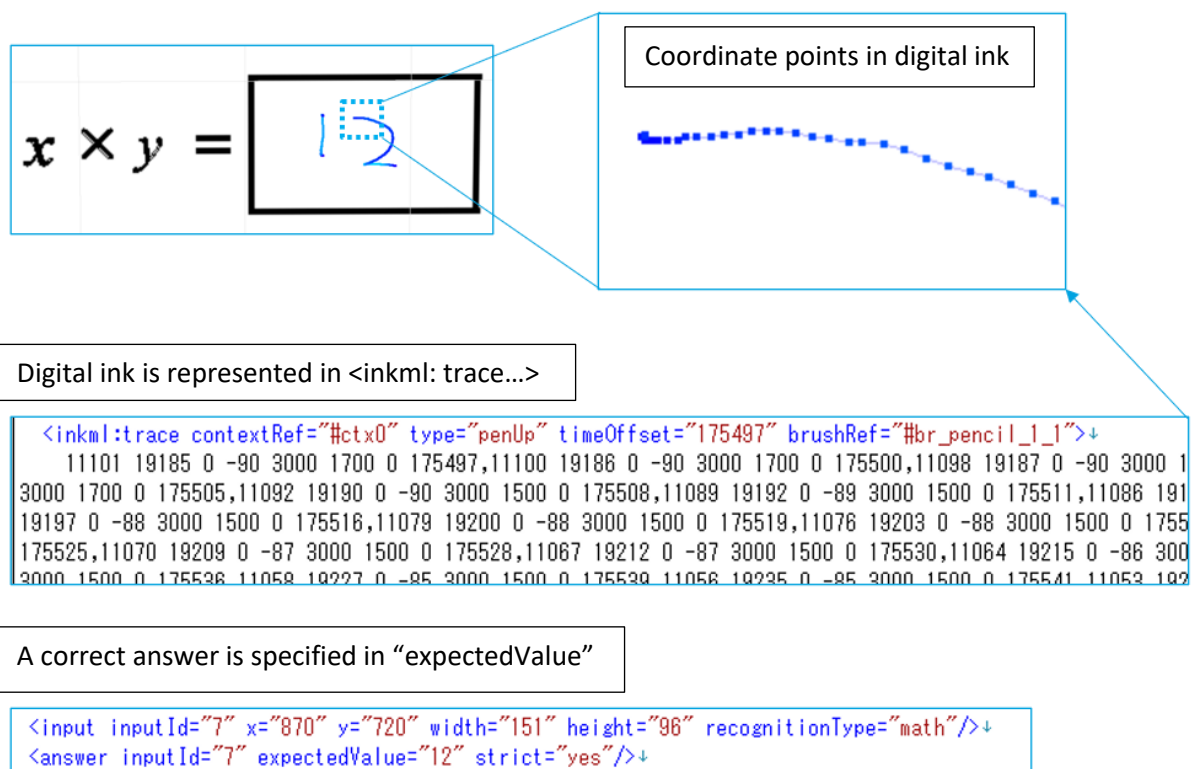**Figure 3:** Digital ink recorded in InkML and correct answer description in XML.

### 3.3.    Issues related to correct answer examples

The workbooks used for this research cover Japanese and mathematics questions in all of grades 1 through 6 of elementary schools and English questions in grades 5 and 6. In Japanese elementary education, English is taught from grade 5. Digitalizing the entire workbooks involves challenges of

processing various types of questions as well as the majority of ordinary questions asking for textual and mathematical answers. We classified all the questions appearing in the workbooks from grades 1 through 6 into 10 types, as listed in **Table 1**. Precisely speaking, each type of question denotes the type of answers expected. Therefore, we use "type of answers." Among them, seven types require line drawing answers and the remaining three types require Japanese text, English text, and math expressions. **Figure 4** shows examples of the connecting and geometry-type answers.

**Table 1**
Ten types of answers.

| Type ID | Type of answer expected |
|---|---|
| Japanese | Japanese text |
| English | English text |
| Math | Mathematical expressions |
| Geometry | Completing a geometric shape by drawing additional lines |
| Connecting | Connecting labels with lines |
| Locating | Locating characters/words/phrases by drawing sidelines or underlines |
| Selecting | Selecting answers by gestures |
| Filling | Filling or painting objects such as circles, squares, and so on |
| Development | Unfolding a 3D object to 2D development |
| Plotting | Plotting graphs and charts |

(a) connecting labels with lines          (b) geometric question
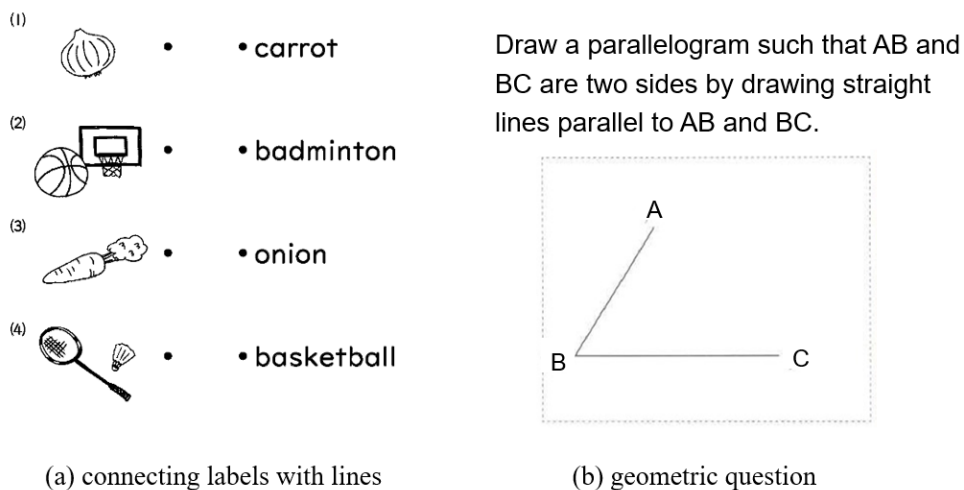
**Figure 4:** Answers of connecting and geometry types (from KUMON publishing with English translation).

A choice format has been used in the workbooks for the questions requiring the selection of a correct stroke order for writing a Kanji character. On the other hand, the electronic paper provides the advantage of obtaining time-series data about handwriting so that the stroke order can be known from the digital ink. If such clear advantages are available, they are chosen, but basically, the questions are presented in the same format as in paper workbooks. Specifically, we followed the styles of paper workbooks as much as possible, but we employed the new styles of questions that exploited the benefits of the electronic paper.

## 3.4.    Data collection

Handwritten answers from 300 elementary school students (50 from each grade) were collected with the permission of the ethics committee of our university. We asked the students to answer the questions

by themselves. Therefore, incorrect answers were included. An automatic scoring system must score correct answers as correct and incorrect answers as incorrect. Therefore, actual incorrect answers are useful as well as correct answers for the study of automatic scoring.

## 4. Prototypes of answer recognizers and scorers

We used deep neural network (DNN) models to build three languages of handwritten answer recognizers: Math recognizer [16, 17], English recognizer [18], and Japanese recognizer [19]. According to the type of input patterns, two groups of handwriting recognizers must be prepared. For online recognizers using pen trajectories as input, the DNN models are composed of multiple stacked Bidirectional Long Short-Term Memory (BLSTM) layers and a Connectionist Temporal Classification layer. For offline recognizers using handwritten images as input, the models consist of an encoder and a decoder with attention layers, where the encoder has multiple convolution layers with pooling layers and the decoder has stacked BLSTM layers. In the previous studies, these models were developed for general usage and evaluated on common handwriting datasets. Thus, they are appropriate to deploy as handwritten answer recognizers.

The Math recognizer achieved a 52.38% expression recognition rate (ink recognition – online recognizer) and a 66.08% expression recognition rate (image recognition – offline recognizer) for the latest CROHME 2019 dataset. The CROHME dataset contains high level mathematical expressions, so much higher recognition rates can be expected in elementary schools. The English text recognizer achieved an 85.78% word recognition rate for IAM-OnDB [20]. The Japanese text recognizer achieved an 86.31% character recognition rate for the TUAT-Kondate dataset [21]. Moreover, the performance of these systems could be improved with transfer learning and dataset adaptation methods, because the above-mentioned recognizers were trained by samples collected in different environments and from different groups of people.

For the other seven types of graphic answers, automatic scorers were prepared by simple methods: dynamic programming matching to compare the answers and expected values for the geometry and plotting types; detecting connecting lines for all the pairs without inconsistency for the connecting type; recognizing drawings inside answer regions and comparing them with the expected values for the selecting and locating types with consistency checking; detection of pen traces for the filling type with consistency checking; and detection of 2D shape objects in specified relations for the development type. The number of samples was too small to prepare DNN scorers. The methods are being elaborated so that a reliable evaluation can be reported later.

Automatic scoring was applied after recognition without any correction of misrecognition or manual labeling. For short answers with one or a few characters, we used perfect matching to provide the scores. For long answers, such as sentences, we employed a scorer that was based on DNNs and the Bidirectional Encoder Representations from Transformers (BERT) model [3].

We performed a preliminary evaluation of automatic Math scoring on the collected dataset of 1st grade elementary school students. The evaluated dataset consisted of 23,848 samples. We employed the precision, recall, and f-measure as the metrics for evaluation, which are defined as follows:

$$Precision = \frac{TP}{TP + FP} \tag{1}$$

$$Recall = \frac{TP}{TP + FN} \tag{2}$$

$$F\text{-}measure = \frac{2}{\frac{1}{Precision} + \frac{1}{Recall}}, \tag{3}$$

where *TP*, *FP*, and *FN* are the numbers of true positive, false positive, and false negative samples, respectively.

**Table 2** shows the results obtained by our online Math recognizer and offline Math recognizer and their combination. We achieved high precision, which means that the system produced a very small *FP*

(incorrect answers that were scored to be correct). On the other hand, we need to increase the Recall by decreasing *FN* (correct answers that were scored to be incorrect). In automatic scoring, *FP* is more serious than *FN*, because the former will not be claimed, but the latter will be claimed by the answerers. The offline recognizer produced inferior performance than the online recognizer but combining them increased Recall to more than 4 percentage point from the best single recognizer.

The Math scorer performed well on the collected dataset, as the handwritten Math answers for elementary grades are much simpler than the handwritten mathematical expressions in CROHME. However, handwriting by children is unstable in stroke order and number, so the combination of online and offline recognizers is effective. In detail, multiple results (candidates) recognized by the online and offline recognizers are re-ranked based on their recognition probabilities. From these candidates, the result with the highest probability is selected as the answer.

**Table 2**

Evaluation of automatic scoring.

|  | Precision (%) | Recall (%) | F-measure (%) |
| --- | --- | --- | --- |
| Online Math recognizer | 99.76 | 89.18 | 94.17 |
| Offline Math recognizer | 99.76 | 82.75 | 90.47 |
| Online + Offline | 99.71 | 94.08 | 96.81 |

## 5. Prototype of user interfaces

**Figure 5** shows the interface that enables the answerers to confirm the scores for their answers and claim any erroneous scores by sending feedback to the system. This feature provides transparency of scoring and encourages answerers to take responsibility for their learning.

We expect that the commitment of the answerers to scoring could be increased, communication with scorers could be promoted to solve problems, and literacy towards AI could be learned, that AI sometimes make mistakes and it is not necessarily perfect, although they must be verified through demonstration experiments.
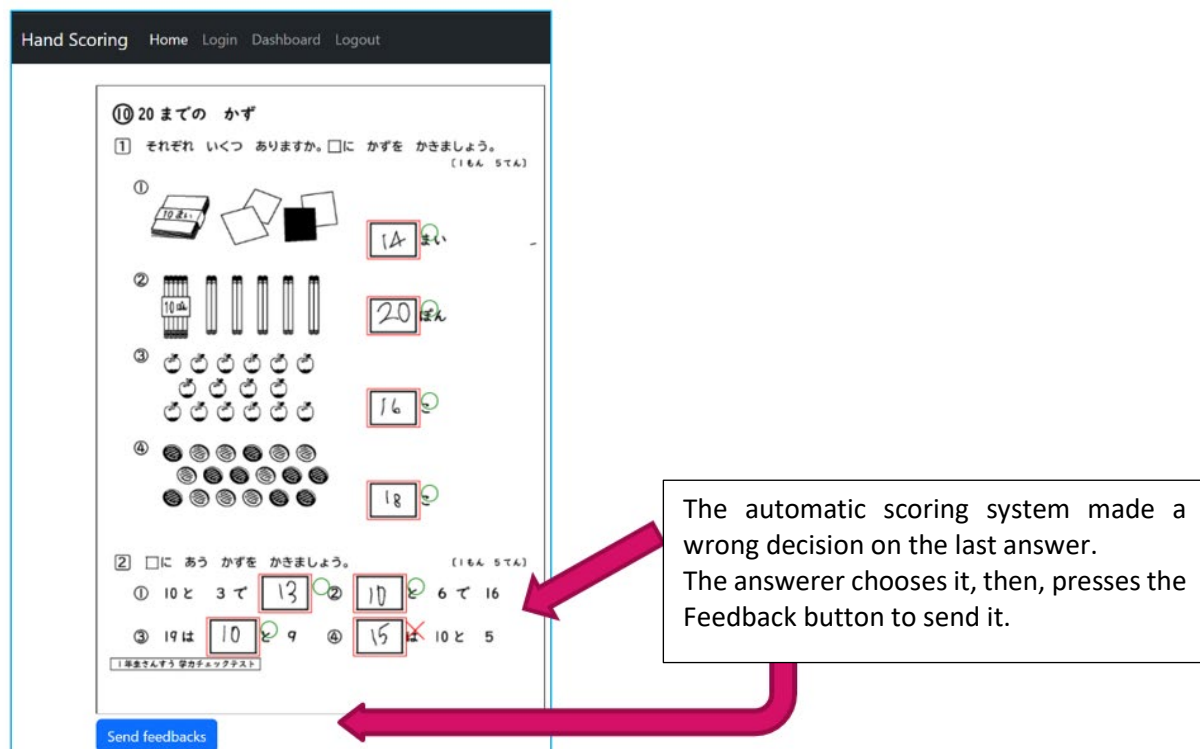


**Figure 5:** Interface of the answerer.

**Figure 6** shows the interface for a human scorer to confirm the scoring results for each question. The answers provided by the answerer are shown in groups of two categories: correct answers bounded in green and incorrect answers bounded in red. This interface highlights the potential of using automated systems to score handwritten answers correctly and efficiently. The human scorer verifies the classification and can correct scores if necessary. Moreover, the human scorer is informed of answers whose scores are claimed to be erroneous by the answerers to confirm the scoring results and take appropriate actions. The scorer revises the score if the answer is scored incorrectly or can help each answerer address the lack of understanding or misunderstanding that led to the incorrect response. This feature helps answerers and scorers create a more collaborative and supportive learning environment.

## 6. Conclusion

We digitalized a set of educational workbooks for primary school on electronic paper and collected digital ink answers from 50 students for every grade in primary school. We also annotated correct answer examples toward automatic scoring of handwritten answers. We combined our existing Japanese, Math, and English recognizers and tailored the automatic scorers. Their performance seems promising. We also created a prototype of the user interfaces for scorers and answerers. In further work, we will perform a feasibility and demonstration experiment for real students and teachers to use the system and show the effect.

The performance of handwritten answer recognizers and scorers should be improved using ensemble recognizers and more powerful DNN models. Interactions between answerers and scorers should be analyzed. This topic is the most important. The user interfaces must be elaborated for real use. This research will be considered successful if the answer scoring can be performed in less time and made less labor-intensive, with higher reliability, and feedback to answerers can be provided quickly for them to review their answers. Even more significant success would be achieved if the answer scoring was
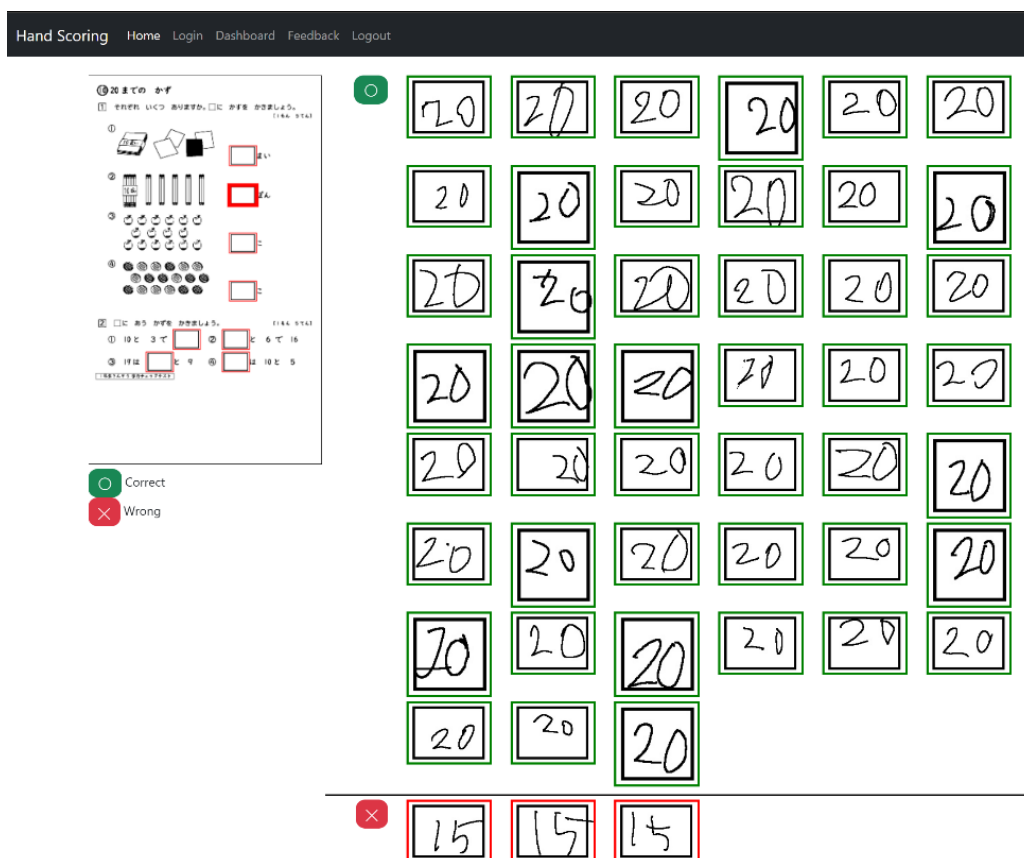


**Figure 6:** Interface of the scorer.

perceived as more transparent, the commitment of the answerers to scoring was increased, communication between answerers and scorers was promoted to solve problems, and AI literacy was learned. The collected ink data has the advantage of revealing the answering process and should be more extensively considered for use.

This research does not contradict the trend of computer-based testing in education, but it will enhance the scope of the questions asked and the role of workbooks to enable answerers to develop creative thinking.

## Acknowledgement

## References

[1] B. Jiang, M. Gu, Y. Du, Recent advances in intelligent textbooks for better learning, in: H. Niemi, R. D. Pea, Y. Lu (Eds.), AI in Learning: Designing the Future, Springer, Cham, 2023, pp. 247-261. doi: 10.1007/978-3-031-09687-7_15.

[2] H. Oka, H. T. Nguyen, C. T. Nguyen, M. Nakagawa, T. Ishioka, Fully automated short answer scoring of the trial tests for common entrance examinations for Japanese university, in: M. M. Rodrigo, N. Matsuda, A. I. Cristea, V. Dimitrova (Eds.), Artificial Intelligent in Education, AIED 2022, Lecture Notes in Computer Science, vol. 13355, Springer, Cham, Durham, UK, 2022, pp. 180-192. doi: 10.1007/978-3-031-11644-5_15.

[3] H. T. Nguyen, C. T. Nguyen, H. Oka, T. Ishioka, M. Nakagawa, Handwriting recognition and automatic scoring for descriptive answers in Japanese language tests, in: Proceedings of the 18th International Conference on Frontiers in Handwriting Recognition, ICFHR 2022, Springer-Verlag (Berlin, Heidelberg), Hyderabad, India, 2022, pp. 274-284. doi: 10.1007/978-3-031-21648-0_19.

[4] X. Liang, S. Sasaki, C. T. Nguyen, M. Nakagawa, Improvement of a computer automated marking system for online handwritten math answers employing machine recognition, IEICE Technical Report, vol. 118, no. 513, PRMU2018-135, 2019, pp. 13-18.

[5] Proceedings of the First International Workshop on Pen-Based Learning Technologies, PLT 2007, IEEE Computer Society, Catania, Italy, 2007. doi: 10.5555/1338440.

[6] J. C. Prey, R. H. Reed, D. A. Berque (Eds.), The impact of tablet PCs and pen-based technology on education 2007: Beyond the Tipping Point, Purdue University Press, 2007.

[7] M. Nakagawa, N. Lozano, H. Oda, Paper architecture and an exam scoring application, in: Proceedings of the First International Workshop on Pen-based Learning Technologies, PLT 2007, Catania, Italy, 2007, pp. 1-6. doi: 10.1109/PLT.2007.15.

[8] N. Lozano, K. Hirosawa, M. Nakagawa, A scoring tool for electronic paper exams, in: Proceedings of the 7th IEEE International Conference on Advanced Learning Technologies, ICALT 2007, Niigata, Japan, 2007, pp. 120-121. doi: 10.1109/ICALT.2007.34.

[9] N. Yoshida, K. Koyama, K. Ng, W. Tsukahara, M. Nakagawa, New features for a pen and paper-based exam scripts marking system, in: T. Bastiaens, J. Dron & C. Xin (Eds.), Proceedings of World Conference on E-Learning in Corporate, Government, Healthcare & Higher Education, E-LEARN 2009, Vancouver, Canada, 2009, pp. 3758-3765.

[10] K. Koyama, M. Nakagawa, Implementation of a pen and paper based exam marking system, in: J. Sanchez, K. Zhang (Eds.), Proceedings of World Conference on E-Learning in Corporate, Government, Healthcare & Higher Education, E-LEARN 2010, Orland, USA, 2010, pp. 1073-1078.

[11] S. Srihari, R. Srihari, P. Babu, H. Srinivasan, On the automatic scoring of handwritten essays, in: Proceedings of the 20th International Joint Conference on Artificial Intelligence, IJCAI'07, Hyderabad, India, 2007, pp. 2880-2884. doi: 10.5555/1625275.1625739.

[12] N. T. Heffernan, C. L. Heffernan, The ASSISTments Ecosystem: Building a platform that brings scientists and teachers together for minimally invasive research on human learning and teaching, International Journal of Artificial Intelligence in Education, 24 (2014), 470-497. doi: 10.1007/s40593-014-0024-x.

[13] U.S. Department of Education, Institute of Education Sciences, What Works Clearinghouse. (2016, June). Secondary Mathematics intervention report: Cognitive Tutor. URL: http://whatworks.ed.gov.

[14] MathNet, 2022. URL: https://www.etrialstestbed.org/projects/mathnet-competition.

[15] W3C, Ink Markup Language (InkML), 2011. URL: https://www.w3.org/TR/InkML/.

[16] C. T. Nguyen, T. N. Truong, H. T. Nguyen, M. Nakagawa, Global context for improving recognition of online handwritten mathematical expressions, in: Proceedings of the 16th International Conference on Document Analysis and Recognition, ICDAR 2021, Lausanne, Switzerland, 2021, pp. 617-631. doi: 10.1007/978-3-030-86331-9_40.

[17] T. N. Truong, C. T. Nguyen, M. Nakagawa, Syntactic data generation for handwritten mathematical expression recognition, Pattern Recognition Letters, 153 (2022) 83-91. doi: 10.1016/j.patrec.2021.12.002.

[18] C. T. Nguyen, M. Nakagawa, Finite state machine based decoding of handwritten text using recurrent neural networks, in: Proceedings of the 15th International Conference on Frontiers in Handwriting Recognition, ICFHR 2016, Shenzhen, China, 2016, pp. 246-251. doi: 10.1109/ICFHR.2016.0055.

[19] H. T. Nguyen, C. T. Nguyen, M. Nakagawa, Online Japanese handwriting recognizers using recurrent neural networks, in: Proceedings of the 16th International Conference on Frontiers in Handwriting Recognition, ICFHR 2018, Niagara Falls, USA, 2018, pp. 435-440. doi: 10.1109/ICFHR-2018.2018.00082.

[20] M. Liwicki, H. Bunke, IAM-OnDB - An on-line English sentence database acquired from handwritten text on a whiteboard, in: Proceedings of the 8th International Conference on Document Analysis and Recognition, ICDAR'05, vol. 2, Seoul, South Korea, 2005, pp. 956-961. doi: 10.1109/ICDAR.2005.132.

[21] T. Matsushita, M. Nakagawa, A database of on-line handwritten mixed objects named "Kondate," in: Proceedings of the 14th International Conference on Frontiers in Handwriting Recognition, ICFHR 2014, Hersonissos, Greece, 2014, pp. 369-374. doi: 10.1109/ICFHR.2014.68.