

On the Effectiveness of Curriculum Learning in Educational Text Scoring

Anonymous EMNLP submission

Abstract

Automatic Text Scoring (ATS) is a widely-investigated task in education. Existing approaches often stressed the structure design of an ATS model and neglected the training process of the model. Considering the difficult nature of this task, we argued that the performance of an ATS model could be potentially boosted by carefully selecting data of varying complexities in the training process. Therefore, we aimed to investigate the effectiveness of *curriculum learning* (CL) in scoring educational text. Specifically, we designed two types of difficulty measurers: (i) *pre-defined*, calculated by measuring a sample’s readability, length, the number of grammatical errors or unique words it contains; and (ii) *automatic*, calculated based on whether a model in a training epoch can accurately score the samples. These measurers were tested in both the *easy-to-hard* to *hard-to-easy* training paradigms. Through extensive evaluations on two widely-used datasets (one for short answer scoring and the other for long essay scoring), we demonstrated that (a) CL indeed could boost the performance of state-of-the-art ATS models, and the maximum improvement could be up to 4.5%, but most improvements were achieved when assessing short and easy answers; (b) the pre-defined measurer calculated based on the number of grammatical errors contained in a text sample tended to outperform the other difficulty measurers across different training paradigms.

1 Introduction

Automatic Text Scoring (ATS) is a common but important task in the field of education. With the aid of ATS techniques, instructors can automatically assess the quality of student-authored texts such as short answers to open-ended questions (Sung et al., 2019a; Lun et al., 2020; Leacock and Chodorow, 2003; Xia et al., 2020; Ramachandran et al., 2015) and relatively longer essays (Uto et al., 2020; Burstein and Chodorow,

1999; Attali and Burstein, 2006; Rodriguez et al., 2019; Amorim et al., 2018). Considering the large student-teacher ratio in certain educational scenarios, e.g., the ratio can be up to 10,000:1 or even worse in Massive Open Online Courses (Pappano, 2012), ATS or writing assessment tools building upon ATS (e.g., AcaWriter (Knight et al., 2020) and Grammarly (Karyuatry, 2018)) have been increasingly used in practice to facilitate instructors’ teaching practices.

Given the important role of ATS in education, researchers have made great efforts in designing effective scoring algorithms with the aid of different techniques, such as early rule-based methods (Leacock and Chodorow, 2003), subsequent machine learning methods with hand-crafted features as input (Mohler et al., 2011; Sultan et al., 2016; Amorim et al., 2018), and recent methods based on deep neural networks that can automatically engineer features from input text (Xia et al., 2020; Rodriguez et al., 2019; Sung et al., 2019b,a; Lun et al., 2020; Uto et al., 2020; Ormerod et al., 2021). In certain writing evaluation tasks, e.g., when assessing students’ responses to the prompt questions, these ATS algorithms have demonstrated scoring performance comparable to human graders. However, there still exist scenarios which call for more research efforts to further improve the performance of ATS (Uto, 2021; Ridley et al., 2020).

It is worth noting that most of the existing ATS studies boost the scoring performance by designing a more dedicated (and oftentimes more complicated) model structure to capture the unique characteristics of a writing assessment task. However, given the difficult nature of this task, i.e., sometimes even experienced human graders can disagree on the score assigned to a piece of writing¹, we

¹For instance, in the short-answer scoring dataset released in the Kaggle platform (<https://www.kaggle.com/competitions/asap-sas/>), there are 12.8% answers which received different scores from two independent experienced human graders; this number is even up to 37.9% in an-

argued that, in addition to designing more dedicated model structure, the process of the model training is worthy of our attention as well and this is where *curriculum learning* (CL) can potentially help. CL is a strategy used to train a prediction model by inputting data sorted in an easy-to-hard manner, which imitates the learning order in human curricula. As a model-agnostic training strategy, CL has been widely investigated and applied in various predictive tasks (Wu et al., 2020; Platanios et al., 2019) in terms of improving a model’s generalization ability and subsequently producing better prediction performance. Inspired by Wang et al. (2021), which indicated that CL can be highly effective in enhancing a supervised prediction model when dealing with a difficult task (e.g., the automatic scoring of student-authored responses), we aimed to investigate how CL can be tapped to improve the performance of ATS in education. Formally, our study was guided by the following **Research Question**:

RQ To what extent can curriculum learning strategies boost the performance of ATS methods used in education?

To answer the above question, we centered our work on the design of the two key components of a CL strategy (Wang et al., 2021; Liu et al., 2018): (i) *difficulty measurer*, which determines the relative difficulty level of a training data sample; and (ii) *training scheduler*, which determines the data subset that should be input to a model in a specific training epoch based on the evaluation from the difficulty measurer. Inspired by previous works on proposing effective CL strategies in the broader NLP research (e.g., Spelling Error Correction (Gan et al., 2021) and Natural Answer Generation (Liu et al., 2018)) as well as the works on automatically characterizing textual data in education, we devised two types of CL strategies in this work, i.e., *pre-defined* and *automatic*, which are grouped according to whether any or both of the two key components described above are pre-defined by human experts or automatically learned in a data-driven fashion. It should be noted that these naming terminologies are in line with those summarized by Wang et al. (2021). Specifically, we studied a total of four pre-defined CL strategies, in which the difficulty level of a piece of written text can

other dataset used for automatic essay scoring (<https://www.kaggle.com/c/asap-aes>), which is essentially a more challenging task.

be measured via calculating its length, readability, the number of grammatical errors or unique words it contains, and the training scheduler is defined as the linear continuous schedulers (Wang et al., 2021). As documented in relevant CL studies in computer vision and NLP, in addition to presenting the training data in an easy-to-hard fashion, sometimes a model can achieve better prediction performance by reverting the training order to hard-to-easy (denoted as *anti-curriculum*). As there lacked prior studies on applying CL to tackle the task of AST and it remained largely unknown which training paradigm would benefit the most, we included both the easy-to-hard and hard-to-easy training paradigms to measure the effectiveness of the four CL strategies described above. Through extensive evaluations on two widely-used educational datasets, i.e., one for Automatic Short Answer Scoring (ASAS)² and the other for Automatic Essay Scoring (AES)³, our work demonstrated that: (i) with the aid of CL, the performance of state-of-the-art ATS models can be further boosted with a maximum 4.5% improvement (measured by Quadratic Weighted Kappa); (ii) among the four investigated pre-defined difficulty measurers, the number of grammatical errors tended to give the most robust performance in measuring sample difficulty; (iii) no significant difference was observed between the pre-defined and automatic CL strategies, or between the easy-to-hard and hard-to-easy training paradigms.

2 Related Work

2.1 Automatic Text Scoring in Education

In education, accurate assessment of textual responses authored by students, along with timely and informative feedback carefully crafted by instructors, is an important task in helping students develop effective writing skills and improve their knowledge level (Dikli, 2010). The completion of this task used to rely on manual efforts heavily. However, manual grading has often suffered from issues such as low precision, high inconsistency, and limited scalability (i.e., being unable to provide timely assessment to a large number of students’ responses) (Fazal et al., 2011; Valenti et al., 2003). As a remedy, ATS has been proposed and widely investigated by researchers to facilitate instructors to perform scoring practices (Alikaniotis et al., 2016;

²<https://www.kaggle.com/competitions/asap-sas/>

³<https://www.kaggle.com/c/asap-aes>

176 Ke and Ng, 2019) and it is often used to assess
177 students' responses to short-answer questions and
178 essay prompts, which are denoted as ASAS and
179 AES, respectively.

180 Broadly speaking, most of the existing ATS can
181 be categorized into two categories (Bonthu et al.,
182 2021; Uto et al., 2020). One is often built upon
183 traditional machine learning techniques e.g., Lin-
184 ear Regression (Nau et al., 2017), Support Vector
185 Machine (Gleize and Grau, 2013), and Random
186 Forests (Ishioka and Kameda, 2017), whose per-
187 formance is heavily dependent on the availability
188 and quality of hand-crafted features such as the
189 number of words contained in an answer (Platanios
190 et al., 2019) and the number of distinct words in
191 the answer (Li et al., 2015). The other is empow-
192 ered by the recent deep learning techniques such
193 as Bi-LSTM (Kim et al., 2018) and BERT (Sung
194 et al., 2019b), which can directly transform the raw
195 text input as embedding-based representations to
196 generate an assessment score without the need of
197 manual feature engineering. For instance, driven
198 by the great success achieved by pre-trained lan-
199 guage models in various NLP tasks, Sung et al.
200 (2019b) proposed to couple BERT (Devlin et al.,
201 2018) with a single classification layer and fine-
202 tuned the whole model on a labeled dataset for
203 ASAS, whose scoring performance was up to 0.91
204 measured by weighted average F1. Though cer-
205 tain methods, e.g., those proposed by Taghipour
206 and Ng (2016); Sung et al. (2019b), have demon-
207 strated human-like scoring performance, there still
208 exist scenarios in which further research efforts are
209 needed to develop more accurate ATS techniques,
210 e.g., the cross-prompt scenario in which data from
211 auxiliary prompts is used to trained ATS models for
212 the target prompt(Uto, 2021; Ridley et al., 2020).

213 Noticeably, the ATS models described in the
214 above studies, especially those based on state-of-
215 the-art deep learning techniques, often treated the
216 design of more dedicated model structures as their
217 main means to boost the performance of ATS, while
218 little attention has been given to the training pro-
219 cess of these well-designed models. As explained
220 before, considering the difficult nature of ATS in
221 education, we argued that it might be worthwhile
222 to apply CL to optimize the training process and
223 assist an ATS model to reach its full potential.

2.2 Curriculum Learning 224

225 Curriculum learning (CL) refers to the strategy
226 used to train a prediction model by imitating the
227 meaningful learning order in human curricula, i.e.,
228 presenting the training samples in an easy-to-hard
229 manner so as to enable the model to first optimize
230 an easier version of the target problem and grad-
231 ually consider harder versions, until solving the
232 full target task of interest (Bengio et al., 2009). As
233 indicated before, most of the existing CL strategies
234 consist of two key components, i.e., *difficulty mea-*
235 *asurer* and *training scheduler* (Wang et al., 2021)
236 (or *scoring function* and *pacing function* in other
237 relevant literature). Depending on whether any (or
238 both) of the two components are designed with the
239 aid of human expertise (or data-driven approaches),
240 a strategy can be classified as either *pre-defined*
241 or *automatic* CL. Take difficulty measurer as an
242 example, Platanios et al. (2019) developed a pre-
243 defined strategy in which the difficulty of input
244 text was determined by using its length as a proxy
245 (i.e., the longer the input text, the higher difficulty
246 it has), while Gan et al. (2021) proposed an auto-
247 matic strategy in which the difficulty was measured
248 by calculating its training loss in a specific epoch.

249 Since its inception Bengio et al. (2009), CL
250 has been demonstrated effective in boosting per-
251 formance of various models in the research of
252 computer vision and NLP (Soviany et al., 2020;
253 Spitzkovsky et al., 2010; Tudor Ionescu et al., 2016;
254 Gan et al., 2021; Platanios et al., 2019; Wei et al.,
255 2016; Liu et al., 2018). For instance, when train-
256 ing on imbalance-distributed image data, the Dy-
257 namic Curriculum Learning framework proposed
258 by Wang et al. (2019) employed a two-level cur-
259 riculum schedulers, which consist of a dynamic
260 sampling scheduler that adjusts the data distribu-
261 tion at each time step and balances the importance
262 between the classification loss and the metric learn-
263 ing loss. In a different vein, when performing the
264 task of natural answer generation, Liu et al. (2018)
265 measured text difficulty from the perspective of
266 Grammar (Stanford Parser score⁴) and trained the
267 model first on the simple and low-quality data and
268 then on the complex and high-quality data to gradu-
269 ally learn to generate reliable answers for questions
270 of different complexity, outperforming the state-of-
271 the-art by an average improvement of about 7.5%
272 in terms of accuracy. More worthy of our attention
273 is that, as indicated by Wang et al. (2021), CL can

⁴<http://nlp.stanford.edu/software/parser-faq.shtml>

274 be particularly useful when dealing with difficult
275 tasks, e.g., those involving the use of higher-order
276 cognitive skills to develop solutions, e.g., the task
277 of assessing student-author responses investigated
278 in our study.

279 It should be pointed out that, though most of
280 the existing CL studies posited that a model’s per-
281 formance can be boosted to the maximum degree
282 by adopting the easy-to-hard learning order, there
283 have been some studies (Zhang et al., 2018; Pi et al.,
284 2016; Braun et al., 2017) which demonstrated that,
285 in certain cases, the model could be trained in an op-
286 posite learning order, i.e., from harder data to easier
287 data (also called anti-curriculum learning (Wang
288 et al., 2021)). For instance, Zhang et al. (2018)
289 demonstrated that the hard-to-easy order, compared
290 to its easy-to-hard counterpart, could lead to better
291 model performance in neural machine translation.

292 Though CL has been demonstrated effective, few
293 studies attempted to investigate its effectiveness in
294 enhancing the assessment of textual responses au-
295 thored by students in education, which is often
296 deemed as a challenging and high-stake task (Gierl
297 et al., 2014; Beseiso and Alzahrani, 2020; Cao
298 et al., 2020), motivating us to design effective CL
299 strategies to further enhance the performance of
300 existing ATS models.

301 3 Methods

302 3.1 Tasks and Datasets

303 Our study was centered around two common writ-
304 ing assessment tasks in education, i.e., ASAS and
305 AES. Generally speaking, as the text length of an
306 essay is often much longer than a short answer,
307 AES is often regarded as a more challenging task
308 than ASAS. The two datasets we used were re-
309 leased by the Hewlett Foundation to spur the devel-
310 opment of novel techniques to tackle the two tasks
311 described above, respectively.

312 **Dataset for ASAS⁵.** The dataset consists of about
313 17,000 answers written by students who were
314 mainly from Grade 10 in the US as responses to
315 10 different prompt questions of subjects includ-
316 ing Science, Biology, English, etc. The number of
317 words contained in an answer ranges from 1 to 344,
318 with an average of 41.7. It is noteworthy that each
319 answer was assessed by two independent human
320 graders. The scores given by the first grader are
321 the ground truth that an ATS model should aim

⁵<https://www.kaggle.com/competitions/asap-sas/>

to predict, while the scores given by the second
grader are only used to measure the agreement be-
tween different human graders. Notice that there
are about 12.8% answers which received different
scores from the two graders.

Dataset for AES⁶ The dataset consists of about
13,000 essays written by students who were from
Grade 7 to Grade 10 in the US as responses to eight
different prompts. Similar to the ASAS dataset,
each answer in this dataset was assessed by at least
two human graders. The difference lies in that the
final score to be predicted by a model was deter-
mined based on the scores provided by all human
graders. Notice that there are 37.9% answers which
received different scores from the human graders.

The descriptive statistics of the two datasets can
be found in the Appendix A.1.

3.2 Models

Recall that our goal was to investigate the effec-
tiveness of CL in boosting the performance of ATS
models. To measure the capabilities of various CL
strategies, we selected state-of-the-art models used
for ASAS and AES as the testbeds in this study, as
described below.

Model for ASAS. In line with previous studies (Xia
et al., 2020; Sung et al., 2019b), given the limited
number of scores that can be assigned to an answer
(e.g., ranging from 0 to 3), we tackled the task of
ASAS as a classification problem. We followed
the approaches developed by (Sung et al., 2019b,a;
Lun et al., 2020), which coupled BERT with a sin-
gle classification layer as the scoring model and
fine-tuned the whole model for each of the prompt
question so as to enable the model to capture the
task-specific characteristics and subsequently opti-
mize the scoring performance.

Model for AES. Here, we adopted the approach de-
veloped by (Uto et al., 2020), i.e., augmenting the
BERT-based grader (i.e., the one used for ASAS de-
scribed above) with human-crafted features as input
to maximize the scoring performance. Particularly,
the features were engineered on the essay level, and
the rationale behind this is that, compared to a short
answer, an essay is often much longer and can have
a more complex structure, which may pose chal-
lenges to the BERT model to derive an accurate
representation for the essay. By adding essay-level
features as input, the BERT-based grader was ex-
pected to model the quality of an essay better. This

⁶<https://www.kaggle.com/competitions/asap-aes/>

approach has been reported to achieve state-of-the-art performance in AES (Ormerod et al., 2021). Similar to ASAS, we built one AES for each of the prompts contained in the dataset.

3.3 Curriculum Learning Design

3.3.1 Difficulty measurer

We investigated two types of difficulty measurers, i.e., *pre-defined* and *automatic*, as describe below.

Pre-defined. A key characteristic of pre-defined strategies lies in that the measurement of a training sample’s difficulty often relies on human expertise (Wang et al., 2021). Inspired by relevant studies on analyzing textual data in education (e.g., those characterizing the utterances generated by instructors in online one-on-one tutoring (Lin et al., 2022a,b) or analyzing students’ posts made in discussion forums (Sha et al., 2021)), we designed a total of four difficulty measures for pre-defined CL strategies, as described below:

- Length, which measures the length of a piece of text as a proxy to its difficulty level (Platanios et al., 2019; Spitkovsky et al., 2010). Here, the longer a response is, the more difficult it is considered to be.
- Distinct-1, which counts the number of unique words contained in a response to measure its difficulty level (Li et al., 2015). Here, the more unique words a response has, the more difficult it is considered to be. We acknowledged that longer text might contain more unique words. Thus, we scaled the number of unique words by the length of the textual response as the final difficulty measurer.
- Readability, which calculates the Flesch Reading Ease score (Farr et al., 1951) of a response to measure its difficulty level. A Flesch Reading Ease score is of the range [0, 100], with 0 representing being extremely difficult to read and 100 being extremely easy to read. Therefore, in this measurer, the lower the readability score, the more difficult the response is.
- Errors, which detects the number of grammatical errors and spelling mistakes contained in a response to measure its difficulty level. Here, the more errors a response contains, the more difficult it is considered to be. Similar to Distinct-1, we noticed that the longer the

text, the more errors it might contain. We therefore also scaled this measure by the text length as the final difficulty measurer.

Automatic. Though pre-defined strategies have been demonstrated effective in various application scenarios, they are often plagued by their strong reliance on human expertise to define an appropriate difficulty measurer and an extensive search for effective combinations of difficulty measurer and training scheduler. Therefore, in addition to the four pre-defined CL strategies described above, we further designed an automatic difficulty measurer to dynamically select data samples based on instance-wise training loss and enable a more flexible training process. Specifically, the automatic difficulty measurer used in this study characterizes data samples as either *easy* and *difficult*, which represents the samples whose ground truth scores are *correctly* or *incorrectly* predicted by a model in a training epoch. Let p_{easy}^t and $p_{difficult}^t$ denote the probabilities of an individual *easy* sample and an individual *difficult* sample to be selected for model training at the current t -th epoch, we define r as the ratio between these two probabilities:

$$r = \frac{p_{difficult}^t}{p_{easy}^t} \quad (1)$$

By choosing different values for r , we can enable the strategy to lay different emphasis on the easy and difficult samples. In particular, we explored two different ways to determine the value for r and consequently two variants of the automatic strategy:

- *Static*, which sets r to the same value across different training epochs. During experiments, r was empirically determined by searching in the range of (0, 5] with an interval of 0.1. When $r < 1$ (i.e., $p_{difficult}^t < p_{easy}^t$), easy samples will be more likely to be selected for training; when $r > 1$, difficult samples will be more likely to be selected for training.
- *Adaptive*, the value of r at the current t -th epoch is based on the number of easy and difficult samples in the previous epoch. We denote the set of easy and difficult samples as E and D , respectively, and formally define:

$$r = \frac{p_{difficult}^t}{p_{easy}^t} = \frac{|E|}{|D|}. \quad (2)$$

Note that such a setting of r ensures that when there are relatively a larger portion of easy (or difficult) samples, the strategy tends to select difficult (or easy) samples more often for the subsequent training.

The sum of the sampling probabilities assigned to all training data should be equal to 1 in both cases for each epoch, e.g., in the variant of Adaptive, it should be $p_{easy}^t * |E| + p_{difficult}^t * |D| = 1$.

3.4 Training scheduler

For both of the pre-defined and automatic difficulty measurers described above, we defined the training scheduler by using a function $\lambda(t)$ to map a training epoch number t to a scalar value $\lambda \in (0, 1]$, i.e., only the λ proportion of the easiest samples should be used to training a model at the t -th epoch. Here, the function is formally defined as:

$$\lambda(t) = \frac{t}{T}, \quad (3)$$

where T denotes the total number of epochs throughout the whole training process and $t \in [1, T]$. Essentially, this is a form of the linear continuous schedulers summarized by Wang et al. (2021). We leave the exploration of other forms of training schedulers (e.g., root function (Platanios et al., 2019) and geometric progression function (Penha and Hauff, 2020)) in our future studies.

As explained before, given the inconsistent findings of CL in various scenarios, we fed the training samples not only in an easy-to-hard order but also in a hard-to-easy order (denoted as Anti-CL) to evaluate the effectiveness of the four pre-defined difficulty measurers. That is, only the λ proportion of the most difficult samples should be used to train a model at the t -th epoch.

3.5 Experimental Setup

Feature engineering for the AES model. Following (Uto et al., 2020), we engineered four types of essay-level features as part of the input to empower the AES model described in Sec.3.2, including length-based features, syntactic features, word-level features, and readability features. Note that each feature was standardized to have a mean of 0 and a standard deviation of 1.0.

Baselines. We implemented two baselines for comparisons: (i) Baseline w/o CL, which refer to the vanilla versions of the selected ASAS and AES models without applying any CL strategies; and (ii)

Random curriculum, in which the proportion of samples used at the t -th training epoch is the same as that of a CL strategy (i.e., as defined in Equation (3)), but the samples were randomly selected from the whole dataset, i.e., being in random difficulty order. This was used to scrutinize whether the observed performance change of a model was caused due to the changing sample size in different epochs (Wu et al., 2020).

Model implementation. We constructed the scoring models based on the pre-trained *bert-base-cased* encoder⁷ implemented by the Python package Transformers⁸. Specifically, to obtain the ASAS scoring model, we simply added a classification layer on top of the *bert-base-cased*. For the AES scoring model, the representation vector output by the *bert-base-cased* encoder and the essay-level features were concatenated to reach an augmented representation vector, which served as the input into a regression layer (linear layer with sigmoid activation) to predict the final score. Different from the ASAS scoring model, the AES model training adopts the mean square error (MSE) loss function, where the scores of the training samples are normalized to $[0, 1]$ (rescaled to the original score range at the prediction stage). All the codes used in this study can be accessed via <https://github.com/AnonymousGitHubLink>.

Model training. For each prompt, we randomly split the data into training, validation, and testing sets in the ratio of 70% : 15% : 15%. When training a model, we set the batch-size as 16 and the number of training epochs as 5. We selected the learning rate from $\{2e-5, 3e-5, 5e-5\}$ and Adam with decoupled weight to optimize the model. Note that the above parameter selections were guided by Devlin et al. (2018). The best combinations of parameters for each prompt were determined based on the model’s performance in the validation set (measured by QWK). Each reported result is a mean over 5 independent runs with the same hyperparameters.

Model evaluation. In line with previous works (Uto et al., 2020; Ormerod et al., 2021), we adopted the metric Quadratic Weighted Kappa (QWK) to measure the agreement between the predicted scores derived by an ATS model and the ground truth scores.

⁷It had 12 layers, with 768 neurons in each hidden layer and the number of attention heads is 12.

⁸<https://github.com/huggingface/transformers>

4 Results

Results on ASAS. Table 1 details the results on the ASAS dataset. We can observe that Random curriculum showed no improvement over Baseline w/o CL on average QWK, implying that simply changing the size of the training set over each epoch (time step) cannot guarantee performance improvement. By comparing the results of Baseline w/o CL with those of different CL strategies, we can make several interesting observations. Firstly, when considering the average QWK over all prompts, we noticed that both pre-defined and automatic measurers delivered better scoring performance. For instance, Readability and Errors measurers outperformed Baseline w/o CL in both of the Curr and Anti-Curr schedulers. Besides, both the two automatic strategies outperformed Baseline w/o CL. Among all these strategies, the Static gained the maximum improvement (4.5%) over Baseline w/o CL. Secondly, when scrutinizing the results in each prompt, interestingly, we found that both the pre-defined and automatic strategies seemed to have been effective in certain prompts (e.g., Prompt 2, 3 and 4). Notice that these prompts contained a much larger fraction of answers that can be difficult to be accurately assessed. For example, the fraction of answers which received different scores from human graders in Prompt 3 was 23.9%. Note that Prompt 3 was also the one in which CL strategies presented the greatest improvements over baseline. This suggests that, CL tended to be effective when dealing with challenging tasks, which is in line with the findings reported by Wang et al. (2021). Thirdly, among the four pre-defined difficulty measurers, Readability tended to be superior to the others, i.e., achieving the best performance when employed with the Anti-Curr scheduler and the second best when employed with the Curr scheduler. Finally, the advantage of the Curr scheduler over the Anti-Curr scheduler was unclear, since these two schedulers had each shown some superiority in specific prompts. This is also supported by Zhang et al. (2018); Pi et al. (2016); Braun et al. (2017); Wang et al. (2021), which demonstrated that the easy-to-hard training order is not necessarily better than the order of hard-to-easy, reminding us that a training sample perceived easy by human might not be as easy for machine learning models and vice versa. To summarize, these results together imply that CL brought performance im-

provement of certain extent in scoring relatively short and easy answers.

Results on AES. Table 2 details the results on the AES dataset. Based on Table 2, we can make several observations similar to the ASAS results presented in Table 1. Firstly, no overall improvement was brought by Random curriculum compared to Baseline w/o CL. This corroborates that it is necessary to consider the difficulty of training samples when applying CL strategies to ATS models. Secondly, both the pre-defined and automatic measurers displayed certain improvements over Baseline w/o CL. For instance, Readability, Errors, and Distinct-1 all outperformed Baseline w/o CL with the Anti-Curr scheduler. As for the automatic strategies, only Static was shown to be superior to Baseline w/o CL. Among all these strategies, the best performance was given by the pre-defined measurer Readability with the Anti-Curr scheduler. Thirdly, among the four proposed measurers, Errors tended to be more robust compared to the others with both the Curr and Anti-Curr schedulers, though there is no significant superiority observed between Curr and Anti-Curr. However, it is worth noting that all the observed improvements are rather limited, i.e., less than 1%. Also, when delving into the results for each prompt, we notice that the proposed strategies tended to be more effective in certain prompts (i.e., Prompt 1, 3, 4, 6). Surprisingly, these prompts have a relatively higher fraction of essays (about 70% on average) which received the same score from human graders, while this fraction is only 48.3% for the rest of the prompts (i.e., Prompt 2, 5, 7, 8). Therefore, prompts 1,3,4 and 6 can be regarded as relatively easier to be assessed. This implies that CL might be rather limited in boosting the ATS performance when dealing with particularly difficult tasks. For instance, the fraction of essays in Prompt 7 and Prompt 8 which received the same score from human graders are only 29% and 28%, respectively. When scrutinizing the results on these two prompts, most of the proposed CL strategies showed no performance improvement over Baseline w/o CL.

5 Discussion and Conclusion

Automatic scoring of student-authored responses, e.g., short answers and essays, is a long-standing task in the field of education. Though various models have been proposed to tackle this task, it remained largely unknown whether the performance

Table 1: Results (QWK) on the ASAS dataset. Results in **bold** indicate being superior to that of Baseline w/o CL. Underlined results indicate being superior to that of Random Difficulty. **Curr** and **Anti-Curr** denote the easy-to-hard and hard-to-easy learning orders, respectively. A higher QWK indicates a better model performance.

Prompt ID		1	2	3	4	5	6	7	8	9	10	Avg.	
Bseline w/o CL		0.803	0.603	0.069	0.678	0.657	<u>0.777</u>	<u>0.631</u>	<u>0.612</u>	0.706	<u>0.755</u>	<u>0.629</u>	
Random difficulty		0.813	0.630	0.115	0.689	0.689	0.696	0.604	0.574	0.715	0.740	0.627	
Pre-defined	Curr	Readability	0.818	0.653	0.052	0.700	0.641	0.789	0.648	0.573	0.720	0.768	0.636
		Length	0.806	0.654	0.149	0.692	0.620	0.614	<u>0.627</u>	<u>0.583</u>	0.718	0.710	0.618
		Errors	0.817	0.618	0.092	0.698	0.732	<u>0.747</u>	<u>0.615</u>	<u>0.597</u>	0.701	0.725	0.634
		Distinct-1	0.799	0.636	0.126	0.718	0.666	0.782	0.641	<u>0.577</u>	0.712	0.710	0.637
	Anti-Curr	Readability	0.786	0.657	0.197	0.701	0.741	0.685	0.594	<u>0.585</u>	0.722	<u>0.741</u>	0.641
		Length	0.766	0.601	0.045	0.706	0.723	<u>0.761</u>	<u>0.616</u>	0.546	0.726	0.694	0.618
		Errors	0.766	0.598	0.167	0.677	0.680	<u>0.708</u>	0.672	<u>0.597</u>	0.696	<u>0.754</u>	0.631
		Distinct-1	0.811	0.677	0.121	0.699	0.628	0.672	0.649	<u>0.576</u>	0.706	<u>0.750</u>	<u>0.629</u>
Automatic	Static	0.826	0.664	0.208	0.698	0.712	0.821	0.633	<u>0.588</u>	0.688	0.730	0.657	
	Adaptive	0.764	0.673	0.182	0.682	0.733	<u>0.776</u>	0.597	<u>0.575</u>	0.691	0.733	0.641	

Table 2: Results (QWK) on the AES dataset. Results in **bold** indicate being superior to that of Baseline w/o CL. Underlined results indicate being superior to that of Random Difficulty. **Curr** and **Anti-Curr** denote the easy-to-hard and hard-to-easy learning orders, respectively. A higher QWK indicates a better model performance.

Prompt ID		1	2	3	4	5	6	7	8	Avg.	
Bseline w/o CL		0.748	<u>0.618</u>	0.628	0.802	<u>0.788</u>	<u>0.793</u>	<u>0.823</u>	<u>0.677</u>	<u>0.735</u>	
Random difficulty		0.770	0.585	0.659	0.813	0.769	0.789	0.771	0.669	0.728	
Pre-defined	Curr.	Readability	0.761	0.627	0.632	0.792	<u>0.775</u>	0.798	<u>0.820</u>	0.617	<u>0.729</u>
		Length	0.765	<u>0.614</u>	0.698	0.815	0.795	0.787	<u>0.818</u>	0.577	<u>0.734</u>
		Errors	0.768	<u>0.616</u>	0.645	0.832	0.769	0.778	0.811	0.670	0.736
		Distinct-1	0.781	0.644	0.639	0.803	<u>0.777</u>	0.775	<u>0.806</u>	0.628	<u>0.731</u>
	Anti-Curr.	Readability	0.742	<u>0.615</u>	0.698	0.820	0.790	0.810	<u>0.812</u>	0.658	0.743
		Length	0.745	0.575	0.673	0.807	0.799	0.804	<u>0.819</u>	0.631	<u>0.732</u>
		Errors	0.788	0.563	0.663	0.806	<u>0.758</u>	0.806	<u>0.818</u>	0.683	0.736
		Distinct-1	0.777	0.579	0.675	0.814	0.780	0.804	0.833	0.638	0.737
Automatic	Static	0.779	0.639	0.685	0.801	0.790	<u>0.790</u>	<u>0.818</u>	0.621	0.740	
	Adaptive	0.800	<u>0.616</u>	0.656	0.793	<u>0.770</u>	0.800	0.824	0.592	<u>0.731</u>	

660 of these models could be further boosted by carefully
661 selecting data in the training process. In this
662 study, we therefore investigated the effectiveness
663 of CL strategies in empowering the performance of
664 ATS in the tasks of ASAS and AES. Specifically,
665 we designed a set of four pre-defined measurers
666 and one automatic measurer to describe the diffi-
667 culty of a data sample, and investigated their effec-
668 tiveness in both the easy-to-hard and hard-to-easy
669 training paradigms. Through extensive evaluations
670 on two different datasets, we showed that CL in-
671 deed can boost the performance of existing state-
672 of-the-art ATS models used in education and the
673 number of grammatical errors contained in a textual
674 response can be used as an effective metric to mea-
675 sure the training difficulty of the response. More
676 importantly, we demonstrated that, when assess-
677 ing relatively short and easy answers, CL tended

678 to display a stronger power in empowering ATS
679 models and the brought improvement can be up
680 to 4.5% measured in QWK. However, when deal-
681 ing with relatively longer and more challenging
682 essays, CL showed little improvement compared to
683 the baselines. To understand the reason behind this
684 and further improve ATS models, one future direc-
685 tion to improve this study is to dissect and observe
686 how those ATS models change during the training
687 process (e.g., the attention weights given to the in-
688 put text across different training epochs), based on
689 which better CL strategies can be developed. In
690 addition, as we only investigated one type of train-
691 ing scheduler in this study, it would be worthwhile
692 to incorporate more advanced scheduling strate-
693 gies (e.g., those proposed in [Platanios et al. \(2019\)](#);
694 [Penha and Hauff \(2020\)](#)) to further empower the
695 ATS models, especially for the task of AES.

696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750

References

- Dimitrios Alikaniotis, Helen Yannakoudakis, and Marek Rei. 2016. Automatic text scoring using neural networks. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 715–725.
- Evelin Amorim, Marcia Cançado, and Adriano Veloso. 2018. Automated essay scoring in the presence of biased ratings. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 229–237.
- Yigal Attali and Jill Burstein. 2006. Automated essay scoring with e-rater® v. 2. *The Journal of Technology, Learning and Assessment*, 4(3).
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48.
- Majdi Beseiso and Saleh Alzahrani. 2020. An empirical analysis of bert embedding for automated essay scoring. *Int. J. Adv. Comput. Sci. Appl.*, 11(10):204–210.
- Sridevi Bonthu, S Rama Sree, and MHM Krishna Prasad. 2021. Automated short answer grading using deep learning: A survey. In *International Cross-Domain Conference for Machine Learning and Knowledge Extraction*, pages 61–78. Springer.
- Stefan Braun, Daniel Neil, and Shih-Chii Liu. 2017. A curriculum learning method for improved noise robustness in automatic speech recognition. In *2017 25th European Signal Processing Conference (EU-SIPCO)*, pages 548–552. IEEE.
- Jill Burstein and Martin Chodorow. 1999. Automated essay scoring for nonnative english speakers. In *Computer mediated language assessment and evaluation in natural language processing*.
- Yue Cao, Hanqi Jin, Xiaojun Wan, and Zhiwei Yu. 2020. Domain-adaptive neural automated essay scoring. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1011–1020.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Semire Dikli. 2010. The nature of automated essay scoring feedback. *Calico Journal*, 28(1):99–134.
- James N Farr, James J Jenkins, and Donald G Paterson. 1951. Simplification of flesch reading ease formula. *Journal of applied psychology*, 35(5):333.
- Anhar Fazal, Tharam Dillon, and Elizabeth Chang. 2011. Noise reduction in essay datasets for automated essay grading. In *OTM Confederated International Conferences "On the Move to Meaningful Internet Systems"*, pages 484–493. Springer.
- Zifa Gan, Hongfei Xu, and Hongying Zan. 2021. Self-supervised curriculum learning for spelling error correction. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3487–3494, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Mark J Gierl, Syed Latifi, Hollis Lai, André-Philippe Boulais, and André De Champlain. 2014. Automated essay scoring and the future of educational assessment in medical education. *Medical education*, 48(10):950–962.
- Martin Gleize and Brigitte Grau. 2013. Limsiiles: Basic english substitution for student answer assessment at semeval 2013. In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, volume 2, pages 598–602.
- Tsunenori Ishioka and Masayuki Kameda. 2017. Over-writable automated japanese short-answer scoring and support system. In *Proceedings of the International Conference on Web Intelligence*, pages 50–56.
- Laksnoria Karyuatry. 2018. Grammarly as a tool to improve students’ writing quality: Free online-proofreader across the boundaries. *JSSH (Jurnal Sains Sosial dan Humaniora)*, 2(1):83–89.
- Zixuan Ke and Vincent Ng. 2019. Automated essay scoring: A survey of the state of the art. In *IJCAI*, volume 19, pages 6300–6308.
- Byung-Hak Kim, Ethan Vizitei, and Varun Ganapathi. 2018. Gritnet: Student performance prediction with deep learning.
- Simon Knight, Antonette Shibani, Sophie Abel, Andrew Gibson, and Philippa Ryan. 2020. Acawriter: A learning analytics tool for formative feedback on academic writing. *Journal of Writing Research*.
- Claudia Leacock and Martin Chodorow. 2003. C-rater: Automated scoring of short-answer questions. *Computers and the Humanities*, 37(4):389–405.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2015. A diversity-promoting objective function for neural conversation models. *arXiv preprint arXiv:1510.03055*.
- Jionghao Lin, Mladen Rakovic, David Lang, Dragan Gasevic, and Guanliang Chen. 2022a. Exploring the politeness of instructional strategies from human-human online tutoring dialogues. In *LAK22: 12th International Learning Analytics and Knowledge Conference*, pages 282–293.
- Jionghao Lin, Shaveen Singh, Lele Sha, Wei Tan, David Lang, Dragan Gašević, and Guanliang Chen. 2022b. Is it a good move? mining effective tutoring strategies from human–human tutorial dialogues. *Future Generation Computer Systems*, 127:194–207.

806	Cao Liu, Shizhu He, Kang Liu, Jun Zhao, et al. 2018.	<i>artificial intelligence in education</i> , pages 381–394.	860
807	Curriculum learning for natural answer generation.	Springer.	861
808	In <i>IJCAI</i> , pages 4223–4229.		
809	Jiaqi Lun, Jia Zhu, Yong Tang, and Min Yang. 2020.	Petru Soviany, Claudiu Ardei, Radu Tudor Ionescu, and	862
810	Multiple data augmentation strategies for improving	Marius Leordeanu. 2020. Image difficulty curricu-	863
811	performance on automatic short answer scoring.	lum for generative adversarial networks (cugan). In	864
812	In <i>Proceedings of the AAAI Conference on Artificial</i>	<i>Proceedings of the IEEE/CVF Winter Conference on</i>	865
813	<i>Intelligence</i> , volume 34, pages 13389–13396.	<i>Applications of Computer Vision</i> , pages 3463–3472.	866
814	Michael Mohler, Razvan Bunescu, and Rada Mihalcea.	Valentin I Spitkovsky, Hiyan Alshawi, and Dan Juraf-	867
815	2011. Learning to grade short answer questions using	sky. 2010. From baby steps to leapfrog: How “less	868
816	semantic similarity measures and dependency graph	is more” in unsupervised dependency parsing. In	869
817	alignments. In <i>Proceedings of the 49th annual meet-</i>	<i>Human Language Technologies: The 2010 Annual</i>	870
818	<i>ing of the association for computational linguistics:</i>	<i>Conference of the North American Chapter of the As-</i>	871
819	<i>Human language technologies</i> , pages 752–762.	<i>sociation for Computational Linguistics</i> , pages 751–	872
		759.	873
820	Jonathan Nau, Aluizio Haendchen Filho, and Guilherme	Md Arafat Sultan, Cristobal Salazar, and Tamara Sum-	874
821	Passero. 2017. Evaluating semantic analysis meth-	ner. 2016. Fast and easy short answer grading with	875
822	ods for short answer grading using linear regression.	high accuracy. In <i>Proceedings of the 2016 Confer-</i>	876
823	<i>Sciences</i> , 3(2):437–450.	<i>ence of the North American Chapter of the Associ-</i>	877
		<i>ation for Computational Linguistics: Human Lan-</i>	878
824	Christopher M Ormerod, Akanksha Malhotra, and Amir	<i>guage Technologies</i> , pages 1070–1075.	879
825	Jafari. 2021. Automated essay scoring using efficient	Chul Sung, Tejas Dhamecha, Swarnadeep Saha, Tengfei	880
826	transformer-based language models. <i>arXiv preprint</i>	Ma, Vinay Reddy, and Rishi Arora. 2019a. Pre-	881
827	<i>arXiv:2102.13136</i> .	training bert on domain resources for short answer	882
828	Laura Pappano. 2012. The year of the mooc. <i>The New</i>	grading. In <i>Proceedings of the 2019 Conference on</i>	883
829	<i>York Times</i> , 2(12):2012.	<i>Empirical Methods in Natural Language Processing</i>	884
		<i>and the 9th International Joint Conference on Natu-</i>	885
830	Gustavo Penha and Claudia Hauff. 2020. Curriculum	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages	886
831	learning strategies for ir. In <i>European Conference on</i>	6071–6075.	887
832	<i>Information Retrieval</i> , pages 699–713. Springer.	Chul Sung, Tejas Indulal Dhamecha, and Nirmal Mukhi.	888
833	Te Pi, Xi Li, Zhongfei Zhang, Deyu Meng, Fei Wu,	2019b. Improving short answer grading using	889
834	Jun Xiao, and Yueting Zhuang. 2016. Self-paced	transformer-based pre-training. In <i>International Con-</i>	890
835	boost learning for classification. In <i>IJCAI</i> , pages	<i>ference on Artificial Intelligence in Education</i> , pages	891
836	1932–1938.	469–481. Springer.	892
837	Emmanouil Antonios Platanios, Otilia Stretcu, Graham	Kaveh Taghipour and Hwee Tou Ng. 2016. A neural	893
838	Neubig, Barnabás Póczos, and Tom M Mitchell. 2019.	approach to automated essay scoring. In <i>Proceed-</i>	894
839	Competence-based curriculum learning for neural	<i>ings of the 2016 conference on empirical methods in</i>	895
840	machine translation. In <i>NAACL-HLT (1)</i> .	<i>natural language processing</i> , pages 1882–1891.	896
841	Lakshmi Ramachandran, Jian Cheng, and Peter Foltz.	Radu Tudor Ionescu, Bogdan Alexe, Marius Leordeanu,	897
842	2015. Identifying patterns for short answer scoring	Marius Popescu, Dim P Papadopoulos, and Vittorio	898
843	using graph-based lexico-semantic text matching. In	Ferrari. 2016. How hard can it be? estimating the	899
844	<i>Proceedings of the Tenth Workshop on Innovative</i>	difficulty of visual search in an image. In <i>Proceed-</i>	900
845	<i>Use of NLP for Building Educational Applications</i> ,	<i>ings of the IEEE Conference on Computer Vision and</i>	901
846	pages 97–106.	<i>Pattern Recognition</i> , pages 2157–2166.	902
847	Robert Ridley, Liang He, Xinyu Dai, Shujian Huang,	Masaki Uto. 2021. A review of deep-neural automated	903
848	and Jiajun Chen. 2020. Prompt agnostic essay	essay scoring models. <i>Behaviormetrika</i> , 48(2):459–	904
849	scorer: A domain generalization approach to cross-	484.	905
850	prompt automated essay scoring. <i>arXiv preprint</i>	Masaki Uto, Yikuan Xie, and Maomi Ueno. 2020.	906
851	<i>arXiv:2008.01441</i> .	Neural automated essay scoring incorporating hand-	907
852	Pedro Uria Rodriguez, Amir Jafari, and Christopher M	crafted features. In <i>Proceedings of the 28th Inter-</i>	908
853	Ormerod. 2019. Language models and automated	<i>national Conference on Computational Linguistics</i> ,	909
854	essay scoring. <i>arXiv preprint arXiv:1909.09482</i> .	pages 6077–6088.	910
855	Lele Sha, Mladen Rakovic, Alexander Whitelock-	Salvatore Valenti, Francesca Neri, and Alessandro Cuc-	911
856	Wainwright, David Carroll, Victoria M Yew, Dra-	chiarelli. 2003. An overview of current research on	912
857	gan Gasevic, and Guanliang Chen. 2021. Assessing	automated essay grading. <i>Journal of Information</i>	913
858	algorithmic fairness in automatic classifiers of edu-	<i>Technology Education: Research</i> , 2(1):319–330.	914
859	cational forum posts. In <i>International conference on</i>		

915 Xin Wang, Yudong Chen, and Wenwu Zhu. 2021. A
 916 survey on curriculum learning. *IEEE Transactions*
 917 *on Pattern Analysis and Machine Intelligence*.

918 Yiru Wang, Weihao Gan, Jie Yang, Wei Wu, and Junjie
 919 Yan. 2019. Dynamic curriculum learning for im-
 920 balanced data classification. In *Proceedings of the*
 921 *IEEE/CVF International Conference on Computer*
 922 *Vision (ICCV)*.

923 Yunchao Wei, Xiaodan Liang, Yunpeng Chen, Xiaohui
 924 Shen, Ming-Ming Cheng, Jiashi Feng, Yao Zhao, and
 925 Shuicheng Yan. 2016. Stc: A simple to complex
 926 framework for weakly-supervised semantic segmen-
 927 tation. *IEEE transactions on pattern analysis and*
 928 *machine intelligence*, 39(11):2314–2320.

929 Xiaoxia Wu, Ethan Dyer, and Behnam Neyshabur.
 930 2020. When do curricula work? *arXiv preprint*
 931 *arXiv:2012.03107*.

932 Linzhong Xia, Mingxiang Guan, Jun Liu, Xuemei Cao,
 933 and Dean Luo. 2020. Attention-based bidirectional
 934 long short-term memory neural network for short
 935 answer scoring. In *International Conference on*
 936 *Machine Learning and Intelligent Communications*,
 937 pages 104–112. Springer.

938 Xuan Zhang, Gaurav Kumar, Huda Khayrallah, Kenton
 939 Murray, Jeremy Gwinnup, Marianna J Martindale,
 940 Paul McNamee, Kevin Duh, and Marine Carpuat.
 941 2018. An empirical exploration of curriculum learn-
 942 ing for neural machine translation. *arXiv preprint*
 943 *arXiv:1811.00739*.

944 A Appendix

945 A.1 Dataset Statistics

Table 3: Statistics of the ASAS dataset. Note that **ELA** is short for **English Language Arts**. The inter-rater agreement was measured as the quadratic weighted kappa between scorer1 and scorer2.

Prompt ID	Subject Area	#Answers	Score Range	Inter-rater Agreement	Average Length
1	Science	1,672	0-3	0.950	47.1
2	Science	1278	0-3	0.900	59.2
3	ELA	1808	0-2	0.681	47.7
4	ELA	1657	0-2	0.683	40.2
5	Biology	1795	0-3	0.962	25.1
6	Biology	1797	0-3	0.952	23.4
7	English	1799	0-2	0.959	41.1
8	English	1799	0-2	0.866	53
9	English	1798	0-2	0.782	49.7
10	Science	1640	0-2	0.887	41.4

Table 4: Statistics of the AES dataset. Note that **PNE** is short for **Persuasive / Narrative / Expository** and **SDR** is short for **Source Dependent Responses**. The inter-rater agreement was measured as the quadratic weighted kappa between scorer1 and scorer2.

Prompt	Essay Type	#Essays	Score Range	Inter-rater Agreement	Average Length
1	PNE	1783	2-12	0.721	350
2	PNE	1800	1-6	0.814	350
3	SDR	1726	0-3	0.769	150
4	SDR	1770	0-3	0.851	150
5	SDR	1805	0-4	0.753	150
6	SDR	1800	0-4	0.776	150
7	PNE	1568	0-30	0.721	250
8	PNE	721	0-60	0.629	650