

Automatic Task Requirements Writing Evaluation via Machine Reading Comprehension

Shiting Xu, Guowei Xu, Peilei Jia, Wenbiao Ding*, Zhongqin Wu, Zitao Liu

TAL Education Group, Beijing, China
{xushiting, xuguowei, jiapeilei, dingwenbiao, wuzhongqin,
liuzitao}@tal.com

Abstract. Task requirements (TRs) writing is an important question type in Key English Test and Preliminary English Test. A TR writing question may include multiple requirements and a high-quality essay must respond to each requirement thoroughly and accurately. However, the limited teacher resources prevent students from getting detailed grading instantly. The majority of existing automatic essay scoring systems focus on giving a holistic score but rarely provide reasons to support it. In this paper, we proposed an end-to-end framework based on machine reading comprehension (MRC) to address this problem to some extent. The framework not only detects whether an essay responds to a requirement question, but clearly marks where the essay answers the question. Our framework consists of three modules: question normalization module, ELECTRA based MRC module and response locating module. We extensively explore state-of-the-art MRC methods. Our approach achieves 0.93 accuracy score and 0.85 F1 score on a real-world educational dataset. To encourage reproducible results, we make our code publicly available at https://github.com/aied2021TRMRC/AIED_2021_TRMRC_code.

Keywords: Task requirements writing · Machine reading comprehension · Pre-training language model · Neural networks.

1 Introduction

Key English Test¹ (KET) and Preliminary English Test² (PET) are examinations to assess the communication ability of the test taker in practical situations. In PET and KET, there are a variety of question types, including speaking, reading, listening, and writing. In writing questions, examinees are not only required to write an essay precisely and correctly but need to make responses to the **Task Requirements (TRs)**. According to official scoring instructions, an essay with poor task achievements should be assigned a low grade. Some examples of TR writing questions are shown in Table 1.

* Corresponding Author: Wenbiao Ding

¹ <https://www.cambridgeenglish.org/exams-and-tests/key>

² <https://www.cambridgeenglish.org/exams-and-tests/preliminary>

Table 1: Examples of TR writing questions in PET.

	<i>A TV company came to your school yesterday to make film.</i>
	<i>Write an email to your English friend Alice. In your email, you should</i>
1	<i>* explain why the TV company chose your school</i>
	<i>* tell her who or what they filmed</i>
	<i>* say when the programme will be shown on television.</i>
	<i>You arranged to meet your English friend Sally next Tuesday,</i>
	<i>but you have to change the time.</i>
2	<i>Write an email to Sally. In your email, you should</i>
	<i>* suggest a new time to meet on Tuesday</i>
	<i>* explain why you need to change the time</i>
	<i>* remind Sally where you arranged to meet</i>

¹ Lines begin with * are task requirement questions.

Timely and accurate evaluation on the performance of test-takers, especially informing them of TR achievements of their essays, is essential to improve their writing and communication skills. Such evaluation usually takes experienced teachers a large amount of time as each essay needs to be graded carefully. However, due to the limitation of teacher resources, most English learners cannot get timely assessments on the quality of their essays. Although many researchers studied how to automatically score an essay, most of the current approaches can only provide total scores without enriched supports [26,6,29]. This is not really helpful for students to improve their writing skills.

In natural language processing field, machine reading comprehension (MRC) has been studied for a long time and can be employed to provide details in terms of how well TRs have been achieved in students' essays. In MRC field, the second version of Stanford Question Answering Dataset (SQuAD 2.0) is the most widely used benchmark dataset to evaluate model performance [22]. However, our experiments prove that even a model that achieves the best performance on SQuAD 2.0 cannot be directly used on educational scenarios, as there is a significant performance degradation. The main reason is that SQuAD 2.0 is a general-purpose open-source dataset, but there is a huge difference between educational and general-purpose corpora.

To alleviate these problems, we construct a real-world educational dataset and propose an end-to-end framework based on MRC approach, which uses ELECTRA as a backbone, to detect whether students respond to TRs in their essays [3]. Our framework can clearly and accurately locate sentences in student essays that respond to the requirements. Experiments on an educational dataset show that the proposed framework achieves 0.93 accuracy score and 0.85 F1 score, outperforming many existing approaches. We believe that this research can help automatic essay scoring system provide interpretable grading results, thereby helping students improve their writing skills.

2 Related Work

2.1 Automated Writing Evaluation

Automated writing evaluation (AWE) has been studied for a long time in both industry and academia [1,20,13,30]. Since Page and Ellis B published their works in 1996, plenty of automated scoring products and applications, e.g., E-rater, have emerged. Based on AWE, lots of works on automatic essay scoring (AES) have been published [20,26,6,29]. However, these works mainly focused on giving a holistic score, which measures the overall quality of an essay. Taghipour explored several neural network models for AWE and outperformed strong baselines without requiring any feature engineering [26]. Dong proposed a reinforcement learning framework that incorporates quadratic weighted kappa as guidance to optimize the scoring system [6]. In recent years, a variety of researches focused on fine-grained essay evaluation [2,12,21]. In Persing’s work, they presented a feature-rich approach to score prompt adherence of essays [21]. In Ke’s work, they not only predicted a score of thesis strength but also provided more reasons [12]. Nevertheless, none of these works address the problem of detecting TR achievements in AES systems.

2.2 Machine Reading Comprehension

At document level, finding students’ response to a TR is similar to extractive and abstractive MRC task in which given several reading materials, the model is expected to answer related questions based on the materials. The MRC models are expected to understand both the context and the question and be able to perform reasoning. In TR writing, we could regard student’s essays as reading materials, and the model is supposed to find answers to TRs. If no answer is found, it indicates that the essay does not respond to the requirement.

The early trend of MRC used long short-term memory or convolutional neural network as an encoder of questions and contexts and blended a variety of attention mechanisms, e.g., attention sum, gated attention [8,19,11,5]. Approaches mimicking the process of how humans do reading comprehension were also proposed, such as multi-hop reasoning [24,16,25]. Recently, pre-trained language models, e.g., BERT, RoBERTa, ALBERT, BART, ELECTRA, became prevalent encoder architectures in MRC and achieved state-of-the-art performance [4,17,14,15,3]. Besides these improvements and optimizations on the encoder module, research about the decoder in the MRC model also starts to draw attention. Zhang et al. proposed an answer verification method and achieved state-of-the-art single model performance on SQUAD 2.0 benchmark with ELECTRA encoder module [31,22,23].

Another line of research on MRC is how to construct high-quality datasets and lots of works have been done [18,7,27,10,23]. Among them, SQuAD is one of the most widely-used reading comprehension benchmarks [23]. However, Rajpurkar et al. showed that the success on SQuAD does not ensure robustness to distracting sentences [22,9]. One reason is that SQuAD focuses on questions for

which a correct answer is guaranteed to exist in the context document. Therefore, models only need to select the span that seems most related to the question, instead of checking that the answer is actually entailed by the text. Based on SQuAD, Rajpurkar et al. proposed SQuAD 2.0. To do well on SQuAD 2.0, systems must not only answer questions when possible but determine when no answer is supported by the paragraph and abstain from answering [22].

Comparing with previous AWE works, to the best of our knowledge, we are the first to use a MRC approach to detect TR achievements in educational domain. We also construct a Student Essay Dataset (SED) which can be deemed as SQuAD 2.0 in the educational field and we explore the usage of a combination of SQuAD 2.0 and SED.

3 Problem Statement

In the TRs writing evaluation task, let Q denote a collection of task requirement questions and q denote a single question in Q . Let t_q^i denotes the i -th token in the question q such that $q = (t_q^1, t_q^2, t_q^3, \dots, t_q^m)$. $E = (t_e^1, t_e^2, t_e^3, \dots, t_e^n)$ is an essay written by a student where t_e^j denotes the j -th token in the essay E . Then the problem is defined as for each requirement q , is there a sequential text span $S = (t_e^j, t_e^{j+1}, \dots, t_e^{j+s})$ in E that responds to the requirement q ? If such span S exists, q is achieved and S needs to be extracted from the essay E , if not, q is not achieved by E .

4 Method

4.1 The Overall Workflow

The overview of our proposed framework is displayed in Figure 1. Our approach is mainly composed of three principal modules, question normalization (QN) module, MRC module, and response locating (RL) module.

4.2 Question Normalization Module

Task requirement questions are proposed from the perspective of examiners, but essays are from examinees’ perspective. This perspective gap brings difficulties to the MRC model. To eliminate the difference, we normalize texts of task requirements with two rule-based methods: switching personal pronouns and deleting redundant words.

Switch Personal Pronouns We use pre-defined rules to replace personal pronouns in the sentence. For example, a question “*What will you do in the summer vacation ?*” may receive a student’s answer “*I will travel to Japan*”. If we change personal pronouns “you” in the question, it will be normalized as “*What will I do in the summer vacation ?*”. The normalized question will decrease the difficulties of this task for the models.

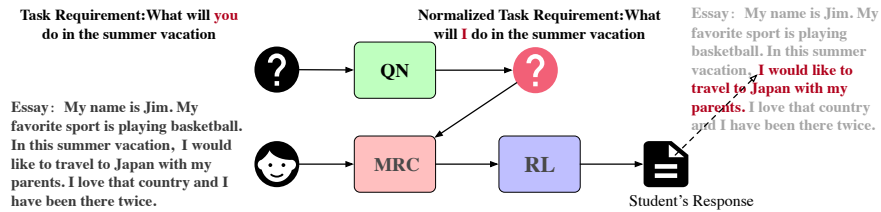


Fig. 1: The workflow of TRs evaluation framework.

Delete Redundant Words We define question words such as “*what*”, “*how*”, etc., and then delete redundant words that appear before them. One example of deleting unnecessary words is that we omit the word “*explain*” in the question “*explain why you need to change the time*” and change it to “*why you need to change the time*”. Another instance is that we delete the words “*remind Sally*” in “*remind Sally where you arranged to meet*” and acquire the normalized question “*Where I arranged to meet*”.

4.3 Machine Reading Comprehension Module

In MRC module, normalized task requirement q and the whole essay E are concatenated with a special symbol $[SEP]$. The entire input sequence to MRC model can be described as $T = ([CLS], t_q^1, t_q^2, t_q^3, \dots, t_q^i, \dots, t_q^m, [SEP], t_e^1, t_e^2, t_e^3, \dots, t_e^j, \dots, t_e^n)$, where the full length of T is $\tau = m + n + 2$.

ELECTRA Encoder We use the discriminator module of ELECTRA to encode each token in T into a dense vector. The max length of T is 512 and tokens exceeding the max length will be truncated at the end. We use h_u^L to represent the final layer outputs of ELECTRA at position u which corresponds to the u -th token in T . We use $H^L = (h_1^L, \dots, h_\tau^L)$ to denote the last-layer hidden states of the input sequence, where $H^L \in \mathbb{R}^{\tau \times 768}$. ELECTRA model is based on a multi-layer bidirectional Transformer encoder, and multi-head attention network [28]. Therefore, h_u^L is able to capture the context of the u -th token from q and E . The attention function in ELECTRA and the output of layer l are showed in eq.(1). In layer l , inputs Q, K, V are computed by $H^{l-1}W_q, H^{l-1}W_k, H^{l-1}W_v$ respectively, where H^{l-1} denotes the output of the previous layer and $W_q \in \mathbb{R}^{768 \times d_k}$, $W_k \in \mathbb{R}^{768 \times d_k}$, $W_v \in \mathbb{R}^{768 \times d_k}$. Thus Q, K, V have the same dimensions $\mathbb{R}^{\tau \times d_k}$ where d_k is the dimension of vectors in K .

$$Attention(Q, K, V) = softmax\left(\frac{Q^T K}{\sqrt{d_k}}\right)V \quad (1)$$

$$H^l = max(0, Attention(Q, K, V)W_1 + b_1)W_2 + b_2$$

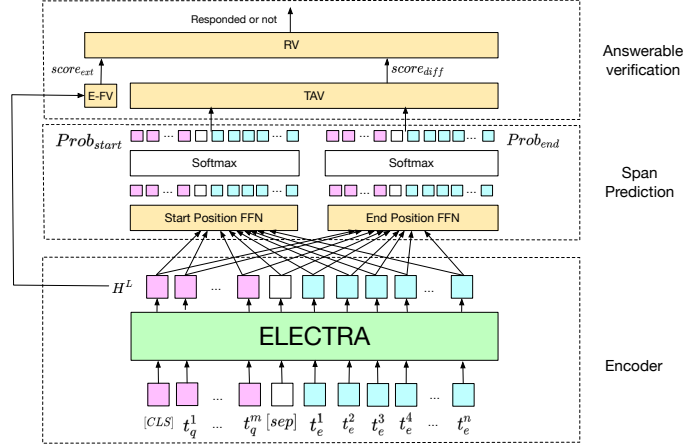


Fig. 2: The architecture of MRC model

Span Prediction We employ a fully connected layer with softmax operation which takes H_L as input and outputs start and end probabilities of each token in T , as shown in eq.(2). Let p_{start}^i and p_{end}^i represent the start and end probabilities of i -th token in T respectively, thus $Prob_{start} = p_{start}^i, i \in [1, \tau]$ is the start probability vector for all tokens in T and $Prob_{end} = p_{end}^i, i \in [1, \tau]$ is the end probability vector for all tokens in T .

$$\begin{aligned} Prob_{start} &= softmax(H^L * W_{start} + b_{start}) \\ Prob_{end} &= softmax(H^L * W_{end} + b_{end}) \end{aligned} \quad (2)$$

Answerable Verification Motivated by Zhang’s work, we introduce the same answerable verification step to determine whether an essay responds to a task requirement [31].

We feed h_1^L which is the representation vector of $[CLS]$ token encoded by ELECTRA into external front verification (E-FV) module. E-FV uses a fully connected layer followed by softmax operation to calculate classification logits $\hat{y}_i = (logit_{ans}, logit_{na})$ where $logit_{ans}$ is a scalar to indicate the answered logits and $logit_{na}$ is a scalar to indicate no-answer logits. We calculate the difference as the external verification score with Equation 3a.

Threshold-based answerable verification (TAV) takes start and end probabilities as input and outputs the no-answer score $score_{diff}$ computed with Equation 3b, 3c and 3d. p_{start}^1 and p_{end}^1 in Equation 3c represents the start and end probabilities of the $[CLS]$ token in T .

Rear verification (RV) combines $score_{diff}$ and $score_{ext}$ to obtain the final answerable score $score_{final}$ as shown in Equation 3e, where β_1 and β_2 are weights.

MRC model predicts that question q is answered by E if $score_{final} > \zeta$, and not answered otherwise, where ζ is a hyper parameter.

$$score_{ext} = \text{logit}_{na} - \text{logit}_{ans} \quad (3a)$$

$$score_{has} = \max(p_{start}^k + p_{end}^l) \quad k, l \in (1, \tau] \text{ and } k \leq l \quad (3b)$$

$$score_{null} = p_{start}^1 + p_{end}^1 \quad (3c)$$

$$score_{diff} = score_{null} - score_{has} \quad (3d)$$

$$score_{final} = \beta_1 score_{diff} + \beta_2 score_{ext}, \quad (3e)$$

4.4 Response Locating Module

In RL module, it takes start probabilities $Prob_{start}$, end probabilities $Prob_{end}$ and answerable score $score_{final}$ as input, and decides the start and end positions according to these inputs. A naive path to achieve this goal is that positions that obtains the highest start and end probabilities are chosen as start and end positions respectively. All tokens between these two positions are extracted as the student’s response to the task requirement. If the start or the end position is less than $m + 1$, in which case a span of question is marked, or their positions are contradictory, e.g., start position greater than end position, the module decides that the question is not responded. Finally, the framework outputs both the binary label indicating whether the student’s essay does respond to the task requirement and the location of the responsive span if it is available.

5 Experiments

5.1 Datasets

SQuAD 2.0 SQuAD 2.0 is the most widely used benchmark in machine reading comprehension literature. It combines the first version of SQuAD with over 50,000 unanswerable questions written adversarially by crowd workers to look similar to answerable ones [23]. It contains 130,319 training examples from 442 Wikipedia articles and 11,873 development examples from 78 Wikipedia articles, where each example is made of a question and an article. This dataset requires that a model should not only answer the question when it is possible but also abstain from answering when there is no answer in the reading materials.

SED This is a real-world student essay dataset that we collect from a third-party K-12 online learning platform. It consists of 9,450 examples in the training set and 3,357 examples in the test set, where each example contains an essay and a requirement question. There are 3,367 different essays and 593 different task requirement questions in the training set. In the test set, the number of essays and requirement questions are 1,655 and 185 respectively. In order to obtain labels, annotators need to firstly decide whether an essay does respond to the

question and label it positive or negative accordingly. Secondly, for all positive essay examples, annotators need to mark the start and end positions of the span in the essay that responds to the question.

Despite that SQuAD 2.0 and SED share similarities in terms of task and structure, there are many differences between them. First of all, SED is in the educational domain and SQuAD 2.0 is from Wikipedia. Secondly, answers in SED are much longer than answers in SQuAD 2.0. Fig 3 illustrates that most answers in SQuAD 2.0 are between 5 to 20 characters, while answers in SED are between 25 to 100 characters. The average length of answers in SQuAD 2.0 is 18.0 while the average length of answers in SED is 103.4. The last difference is that there are more grammatical errors in SED because essays in SED are written by second language learners. So a model that achieves the best performance on SQuAD 2.0 may not be directly deployed on educational scenarios.

5.2 Experimental Setting

In this section, we describe three sets of experiments as follows.

- Set 1. This set aims to prove that existing SOTA models on SQuAD 2.0 cannot be directly deployed on educational scenarios. In Set 1, all models are trained on SQuAD 2.0 but evaluated on the test set of SED. SAN was trained 50 epochs with learning rate $2e^{-3}$ on SQuAD 2.0 [16]. Pre-trained language models such as BERT, RoBERTa, ALBERT, and BART, were acquired from hugging face³. Our ELECTRA-based approach was trained 2 epochs with default parameters in this work [31].
- Set 2. This set is to prove that MRC approaches are effective solutions to TRs writing evaluation when trained on the educational corpus. The training parameters of the models are consistent with those in Set 1. The difference is that models are all trained on SED.
- Set 3. This set explores how can we utilize SQuAD 2.0 and further improve model performance on SED. Following the idea that models pre-trained on massive data can be a good warm-up for subsequent finetuning, we first train MRC models on SQuAD 2.0 so as to acquire basic models, and then finetune them on SED for optimal performance.

³ <https://huggingface.co>

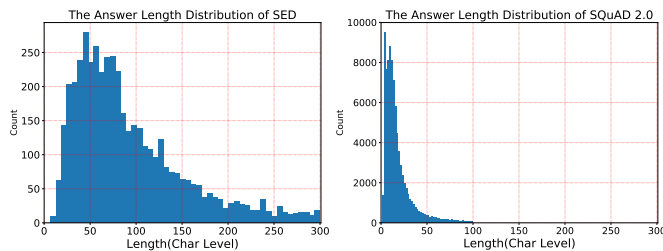


Fig. 3: Distributions of answer length (char level) in SED and SQuAD 2.0.

$$\begin{aligned}
 Accuracy &= \frac{N_{correct}}{N_{total}} \\
 precision &= \frac{Num_{overlap}}{Num_{predict}} \\
 recall &= \frac{Num_{overlap}}{Num_{gold}} \\
 F1 &= \frac{2 * precision * recall}{precision + recall}
 \end{aligned} \tag{4}$$

In all experiments, we use two evaluation indicators. One is Accuracy (Acc.) which measures the performance of the model on the binary classification task of predicting whether the essay answers the TR. Another is Answer Overlap F1 score (F1) which measures the performance of the model to predict the location of the answer span. Accuracy and F1 metrics can be calculated by Equations 4.

In Equation 4, N_{total} indicates the number of examples in the test set, and $N_{correct}$ is the number of examples that are correctly predicted by the framework. $Num_{overlap}$ represents the number of identical tokens in both the predicted span and the gold span. $Num_{predict}$ is the total number of tokens in predicted span and Num_{gold} is the total number of tokens in the gold span.

5.3 Results and Evaluation

Results of Set 1. Table 2 shows that existing SOTA models on SQuAD 2.0 are suffered a significant performance degradation on SED. All models in Table 2 are well finetuned on SQuAD 2.0 and their F1 scores on SQuAD 2.0 dev set are all over 0.66. However, when evaluating them on SED test set, performances drop dramatically. For example, RoBERTa and our method achieve F1 score of 0.83 and 0.89 on SQuAD 2.0 dev set, but both drop to F1 score of 0.49 on SED test set.

Result of Sets 2&3. Table 3 shows results of Set 2 and Set 3. From the results of Set 2, we conclude that the MRC approaches can solve the TRs writing evaluation problem. Comparing with models trained on SQuAD 2.0 (Set 1), models trained on SED achieve significantly better results on SED test set. Our framework in Set 2 achieves the best F1 score of 0.84 and the best accuracy of 0.91, outperforms our framework in Set 1 by 23% Accuracy and 35% F1 score.

Table 2: Performances of models trained on SQuAD 2.0

Training Dataset	Methods	SQuAD 2.0 dev		SED test	
		ACC.	F1	ACC.	F1
SQuAD 2.0 (Set 1)	SAN	0.70	0.66	0.57	0.31
	BERT	0.78	0.73	0.65	0.37
	ALBERT	0.85	0.81	0.58	0.37
	RoBERTa	0.86	0.83	0.69	0.49
	BART	0.87	0.81	0.62	0.42
	Ours	0.92	0.89	0.68	0.49

If we compare results in Set 2 and Set 3, we find that optimal performance can be obtained by firstly training models on SQuAD 2.0 and then finetuning on SED. Specifically, F1 score of SAN increases by 11%, and F1 score of BERT increases by 8%. Similarly, the accuracy also increases significantly in Set 3.

Table 3: Performances of models trained on SED and SQuAD 2.0&SED

Training Dataset	Methods	SED Test	
		Acc.	F1
SED (Set 2)	SAN	0.67	0.58
	BERT	0.79	0.68
	ALBERT	0.84	0.77
	RoBERTa	0.81	0.71
	BART	0.82	0.73
	Ours	0.91	0.84
SQuAD 2.0&SED (Set 3)	SAN	0.79 (+0.12)	0.69 (+0.11)
	BERT	0.84 (+0.05)	0.76 (+0.08)
	ALBERT	0.86 (+0.02)	0.80 (+0.03)
	RoBERTa	0.88 (+0.07)	0.80 (+0.09)
	BART	0.89 (+0.07)	0.82 (+0.09)
	Ours	0.93 (+0.02)	0.85 (+0.01)

Comparing with Set 2, the accuracy of BART and our framework increase by 7% and 2% respectively. Furthermore, our approach achieves the best performance in each of the three sets of experiments, and outperforms a variety of SOTA approaches.

6 Conclusion

In this paper, we proposed a MRC based approach which cannot only detect if an essay responds to a requirement question but find where the essay answers the question. From our experiments and analysis, we demonstrate that SQuAD 2.0 is very different from our educational dataset, so existing SOTA models on SQuAD 2.0 cannot be directly deployed on educational scenarios. Instead, we propose to firstly train a basic model on SQuAD 2.0 and then finetune the basic model on

educational data for optimal performance. We believe this proposed framework is able to help automatic essay scoring systems provide detailed grading results, thereby helping students improve their writing skills.

Acknowledgment

This work was supported in part by National Key R&D Program of China, under Grant No. 2020AAA0104500 and in part by Beijing Nova Program (Z201100006820068) from Beijing Municipal Science & Technology Commission.

References

1. Beigman Klebanov, B., Madnani, N.: Automated evaluation of writing – 50 years and counting. In: Proc. of ACL. pp. 7796–7810. Association for Computational Linguistics (2020). <https://doi.org/10.18653/v1/2020.acl-main.697>, <https://www.aclweb.org/anthology/2020.acl-main.697>
2. Carlile, W., Gurrupadi, N., Ke, Z., Ng, V.: Give me more feedback: Annotating argument persuasiveness and related attributes in student essays. In: Proc. of ACL. pp. 621–631. Association for Computational Linguistics (2018). <https://doi.org/10.18653/v1/P18-1058>, <https://www.aclweb.org/anthology/P18-1058>
3. Clark, K., Luong, M., Le, Q.V., Manning, C.D.: ELECTRA: pre-training text encoders as discriminators rather than generators. In: Proc. of ICLR. OpenReview.net (2020), <https://openreview.net/forum?id=r1xMH1BtvB>
4. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: Proc. of NAACL-HLT. pp. 4171–4186. Association for Computational Linguistics (2019). <https://doi.org/10.18653/v1/N19-1423>, <https://www.aclweb.org/anthology/N19-1423>
5. Dhingra, B., Liu, H., Yang, Z., Cohen, W., Salakhutdinov, R.: Gated-attention readers for text comprehension. In: Proc. of ACL. pp. 1832–1846. Association for Computational Linguistics (2017). <https://doi.org/10.18653/v1/P17-1168>, <https://www.aclweb.org/anthology/P17-1168>
6. Dong, F., Zhang, Y., Yang, J.: Attention-based recurrent convolutional neural network for automatic essay scoring. In: Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017). pp. 153–162. Association for Computational Linguistics (2017). <https://doi.org/10.18653/v1/K17-1017>, <https://www.aclweb.org/anthology/K17-1017>
7. Hewlett, D., Lacoste, A., Jones, L., Polosukhin, I., Fandrianto, A., Han, J., Kelcey, M., Berthelot, D.: WikiReading: A novel large-scale language understanding task over Wikipedia. In: Proc. of ACL. pp. 1535–1545. Association for Computational Linguistics (2016). <https://doi.org/10.18653/v1/P16-1145>, <https://www.aclweb.org/anthology/P16-1145>
8. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural computation* (8) (1997)
9. Jia, R., Liang, P.: Adversarial examples for evaluating reading comprehension systems. In: Proc. of EMNLP. pp. 2021–2031. Association for Computational Linguistics (2017). <https://doi.org/10.18653/v1/D17-1215>, <https://www.aclweb.org/anthology/D17-1215>

10. Joshi, M., Choi, E., Weld, D., Zettlemoyer, L.: TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In: Proc. of ACL. pp. 1601–1611. Association for Computational Linguistics (2017). <https://doi.org/10.18653/v1/P17-1147>, <https://www.aclweb.org/anthology/P17-1147>
11. Kadlec, R., Schmid, M., Bajgar, O., Kleindienst, J.: Text understanding with the attention sum reader network. In: Proc. of ACL. pp. 908–918. Association for Computational Linguistics (2016). <https://doi.org/10.18653/v1/P16-1086>, <https://www.aclweb.org/anthology/P16-1086>
12. Ke, Z., Inamdar, H., Lin, H., Ng, V.: Give me more feedback II: Annotating thesis strength and related attributes in student essays. In: Proc. of ACL. pp. 3994–4004. Association for Computational Linguistics (2019). <https://doi.org/10.18653/v1/P19-1390>, <https://www.aclweb.org/anthology/P19-1390>
13. Ke, Z., Ng, V.: Automated essay scoring: A survey of the state of the art. In: Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019. pp. 6300–6308. ijcai (2019). <https://doi.org/10.24963/ijcai.2019/879>, <https://doi.org/10.24963/ijcai.2019/879>
14. Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., Soricut, R.: ALBERT: A lite BERT for self-supervised learning of language representations. In: Proc. of ICLR. OpenReview.net (2020), <https://openreview.net/forum?id=H1eA7AEtvS>
15. Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., Zettlemoyer, L.: BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In: Proc. of ACL. pp. 7871–7880. Association for Computational Linguistics (2020). <https://doi.org/10.18653/v1/2020.acl-main.703>, <https://www.aclweb.org/anthology/2020.acl-main.703>
16. Liu, X., Shen, Y., Duh, K., Gao, J.: Stochastic answer networks for machine reading comprehension. In: Proc. of ACL. pp. 1694–1704. Association for Computational Linguistics (2018). <https://doi.org/10.18653/v1/P18-1157>, <https://www.aclweb.org/anthology/P18-1157>
17. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692 (2019)
18. Nguyen, T., Rosenberg, M., Song, X., Gao, J., Tiwary, S., Majumder, R., Deng, L., MARCO, M.: A human generated machine reading comprehension dataset. arXiv preprint ArXiv:1607.06275 (2016)
19. O’Shea, K., Nash, R.: An introduction to convolutional neural networks. arXiv preprint arXiv:1511.08458 (2015)
20. Page, E.B.: The imminence of... grading essays by computer. *The Phi Delta Kappan* (5) (1966)
21. Persing, I., Ng, V.: Modeling prompt adherence in student essays. In: Proc. of ACL. pp. 1534–1543. Association for Computational Linguistics (2014). <https://doi.org/10.3115/v1/P14-1144>, <https://www.aclweb.org/anthology/P14-1144>
22. Rajpurkar, P., Jia, R., Liang, P.: Know what you don’t know: Unanswerable questions for SQuAD. In: Proc. of ACL. pp. 784–789. Association for Computational Linguistics (2018). <https://doi.org/10.18653/v1/P18-2124>, <https://www.aclweb.org/anthology/P18-2124>

23. Rajpurkar, P., Zhang, J., Lopyrev, K., Liang, P.: SQuAD: 100,000+ questions for machine comprehension of text. In: Proc. of EMNLP. pp. 2383–2392. Association for Computational Linguistics (2016). <https://doi.org/10.18653/v1/D16-1264>, <https://www.aclweb.org/anthology/D16-1264>
24. Shen, Y., Huang, P., Gao, J., Chen, W.: Reasonet: Learning to stop reading in machine comprehension. In: Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13 - 17, 2017. pp. 1047–1055. ACM (2017). <https://doi.org/10.1145/3097983.3098177>, <https://doi.org/10.1145/3097983.3098177>
25. Shen, Y., Liu, X., Duh, K., Gao, J.: An empirical analysis of multiple-turn reasoning strategies in reading comprehension tasks. In: Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 957–966. Asian Federation of Natural Language Processing (2017), <https://www.aclweb.org/anthology/I17-1096>
26. Taghipour, K., Ng, H.T.: A neural approach to automated essay scoring. In: Proc. of EMNLP. pp. 1882–1891. Association for Computational Linguistics (2016). <https://doi.org/10.18653/v1/D16-1193>, <https://www.aclweb.org/anthology/D16-1193>
27. Trischler, A., Wang, T., Yuan, X., Harris, J., Sordoni, A., Bachman, P., Suleman, K.: NewsQA: A machine comprehension dataset. In: Proceedings of the 2nd Workshop on Representation Learning for NLP. pp. 191–200. Association for Computational Linguistics (2017). <https://doi.org/10.18653/v1/W17-2623>, <https://www.aclweb.org/anthology/W17-2623>
28. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA. pp. 5998–6008 (2017), <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>
29. Wang, Y., Wei, Z., Zhou, Y., Huang, X.: Automatic essay scoring incorporating rating schema via reinforcement learning. In: Proc. of EMNLP. pp. 791–797. Association for Computational Linguistics (2018). <https://doi.org/10.18653/v1/D18-1090>, <https://www.aclweb.org/anthology/D18-1090>
30. Wang, Z., Liu, H., Tang, J., Yang, S., Huang, G.Y., Liu, Z.: Learning multi-level dependencies for robust word recognition. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 9250–9257 (2020)
31. Zhang, Z., Yang, J., Zhao, H.: Retrospective reader for machine reading comprehension. arXiv preprint arXiv:2001.09694 (2020)