# Matching Exemplar as Next Sentence Prediction (MeNSP): Zero-shot Prompt Learning for Automatic Scoring in Science Education

Xuansheng Wu $^1,$  Xinyu He $^2,$  Tianming Liu $^1,$  Ninghao Liu $^1,$  and Xiaoming Zhai $^{2\star}$ 

<sup>1</sup>School of Computing, University of Georgia, Athens, GA, USA <sup>2</sup>AI4STEM Education Center, University of Georgia, Athens, GA, USA {xuansheng.wu, xinyu.he1, tliu, ninghao.liu, xiaoming.zhai}@uga.edu

Abstract. Developing natural language processing (NLP) models to automatically score students' written responses to science problems is critical for science education. However, collecting sufficient student responses and labeling them for training or fine-tuning NLP models is time and cost-consuming. Recent studies suggest that large-scale pre-trained language models (PLMs) can be adapted to downstream tasks without fine-tuning by using prompts. However, no research has employed such a prompt approach in science education. As students' written responses are presented with natural language, aligning the scoring procedure as the next sentence prediction task using prompts can skip the costly finetuning stage. In this study, we developed a zero-shot approach to automatically score student responses via Matching exemplars as Next Sentence Prediction (MeNSP). This approach employs no training samples. We first apply MeNSP in scoring three assessment tasks of scientific argumentation and found machine-human scoring agreements, Cohen's Kappa ranges from 0.30 to 0.57, and F1 score ranges from 0.54 to 0.81. To improve scoring performance, we extend our research to the few-shots setting, either randomly selecting labeled student responses at each grading level or manually constructing responses to fine-tune the models. We find that one task's performance is improved with more samples, Cohen's Kappa from 0.30 to 0.38, and F1 score from 0.54 to 0.59; for the two other tasks, scoring performance is not improved. We also find that randomly selected few-shots perform better than the human expert-crafted approach. This study suggests that MeNSP can yield referable automatic scoring for student-written responses while significantly reducing the cost of model training. This method can benefit low-stakes classroom assessment practices in science education. Future research should further explore the applicability of the MeNSP in different types of assessment tasks in science education and further improve the model performance. Our code is available at https://github.com/JacksonWuxs/MeNSP.

**Keywords:** Prompt Learning, Pre-trained Language Model, written Response, Automatic Scoring, Natural Language Processing.

<sup>\*</sup> Corresponding Author: xiaoming.zhai@uga.edu, ninghao.liu@uga.edu, tliu@uga.edu

#### 1 Introduction

To engage students in meaningful science learning, it is crucial that science assessments can elicit student knowledge-in-use so that instructors can better understand and support student learning [23]. The need for assessing student competence in using scientific knowledge to solve problems or design solutions has been identified in the actualization of  $21^{\text{st}}$ -century science education [4]. While such science and engineering practices are desirable for students to enact, it has been challenging to assess with multiple-choice questions. More complicated written response assessments are needed to assess such scientific practices.

Although the written response assessment allows for the freedom of expression elicited in responses, it accounts for a number of challenges in the process of scoring. Since scoring is expected to be fair and reliable, it requires, among other efforts, that educators are trained using pre-designed rubrics to enhance reliability. Despite the expertise and rubric familiarity of the raters, it demands a lot of time to score the learners' responses appropriately.

In contrast to some of the issues associated with human scoring, machine learning, such as Natural Language Processing (NLP), offers quicker, less expensive, and more consistent scoring [35]. However, most NLP scoring models need to be tested with human scoring as benchmarks before applying in classroom settings. Their validity and performance are a function of the former [33,37]. The majority of the models needed to achieve this require a large sample of data for training, validation, and testing purposes [39]. Given that the conditions necessary to collect a reasonable size of data are a huge effort in the field of education, satisfying the needs of these models is cumbersome [38].

To overcome these challenges, this study employs prompt learning with pretrained language models (PLMs) to develop algorithmic models for automatic scoring. Prompt learning utilizes little or limited labeled data (i.e., zero-shot) in developing a model and has shown great potential in accomplishing NLP tasks with significant efficiency [41]. To verify the efficiency of this new approach, we develop NLP scoring models to automatically score students' written arguments to science phenomena when engaging in scientific argumentation practices.

This study will contribute to the efficiency of developing innovative assessments in education broadly and science education specifically. Our approach will significantly save the time and cost of developing NLP scoring algorithmic models for automatic scoring, which is a bottleneck that has prevented the broad use of machine learning-based assessment practices in classrooms [39]. Implementing this approach could benefit millions of students using automatically scored written response assessments that have been developed [8].

## 2 Related Work

#### 2.1 Natural Language Processing for Automatic Scoring

Natural language processing (NLP) is a field of computer science using computational techniques to learn, understand, and produce human language content.

There has been a long time to apply NLP in language education, such as to correct students' writing errors, conduct semantic analysis, or assess language skills [12]. With the improvement of using NLP to evaluate students' writing skills, research has begun to explore using NLP to evaluate the content of students' writing and study the quality of writing at the level of domain-specifics such as science learning, written response assessment in science education is one of the specific domains needing NLP for automatic scoring [39].

To help teachers better understand students' scientific thinking, researchers have explored using NLP technologies to score student-written responses a decade ago automatically. Haudek et al. [9] used SPSS Text Analytics for Surveys to score students' biology understanding automatically. The program can extract key linguistic features of student writing to classify student responses according to scoring rubrics. Researchers [20] at Carnegie Mellon University developed SIDE (current version named LightSIDE) that integrates various algorithmic functions. Nehm et al. [20] employed his package to develop a portal Evo-Grader to examine students' understanding of biology concepts. Educational Testing Service developed a C-rater [13] for the automatic scoring of GRE essays and short written-response answers. Later, they incorporated and upgraded the tool to C-rater-ML [7,11] and employed it to automatically score students' scientific argumentation and explanations. While prior tools employed individual algorithms, researchers [18,31] also developed tools that ensemble multiple algorithms to score student-written responses simultaneously. Most recently, researchers [25,1] also employed BERT [5] for automatic scoring of student-written responses for scientific practices.

Among these developments, recent surveys [38,39] suggested that developing these models requires a large number of human-scored written responses with varying accuracy. However, collecting and validating the student responses and rigorous rubrics for scoring takes much effort. Since training human experts to use the scoring rubric to assign scores to student responses reliably is challenging, obtaining datasets to develop these models is costly.

#### 2.2 Prompt Learning

Prompt learning [3,30,32,34] leads to a new paradigm in NLP as it can achieve comparable performance to full-parameter fine-tuning with fewer training samples and parameters. "Prompt" typically is a short piece of text that include instructions for the task (zero-shot learning) or a few samples of the task (few-shots learning) [19]. By selecting appropriate prompts, it is possible to directly predict the target label using the pre-trained language models. Compared with the previous "pre-train and fine-tune" paradigm, prompt learning has the advantage of reducing the training cost and being applicable to multiple tasks by changing the prompts [14].

The new paradigm brings its advantage as well as a new challenge, that is, how to find the most appropriate prompt. The most natural way is to create intuitive prompts manually [27,15]. However, some researchers [29,14] found

that the optimal prompt may not be readable by humans. Thus, many methods [29,28,15,16] are proposed to learn a better prompt automatically. Besides learning prompts, providing and ordering a few additional answers in the prompts can also result in satisfied model performance [6,17].

To the best of our knowledge, the application of prompt learning in the education field is still at the beginning. Hart-Davidson et al. (2021) [21] applied prompt learning in the qualitative coding research task, which could be used to provide written feedback for student writing. Zhang et al. (2022) [40] applied prompt learning in the fine-tuning process to boost the performance for automatic scoring of short Math answers. How to give full play to The advantage of using prompt learning hasn't been fully discussed.

In summary, there is a lack of research exploring PLMs' zero-shot performance for automatic scoring in education, not even in science education. Therefore, this study applies prompt learning to automatically score student-written responses to science assessments. The approach can be applied to more tasks in education to save time and cost.

# 3 Approach

This section proposes a method to develop the system for scoring student-written arguments without fully labeled training datasets. Inherently, student responses are presented with natural language, which is readable by PLMs. Since PLMs are built under some pre-trained tasks over natural language corpus, we can reformat the scoring task as one of the pre-trained tasks so that the textual responses of students can be graded without further fine-tuning of the models on labeled responses. Reformatting a downstream task as a pre-trained task is also known as prompt learning, which has been widely applied to handle NLP tasks under the zero-shot and few-shots settings where labeled data are unavailable or limited. However, letting PLMs understand (1) the meaning of the "scoring" procedure, and (2) the scoring rubrics of each assessment item are two non-trivial challenges to achieving the goal. Technically, we need to find a pre-trained task to reformat the scoring procedure and combine student responses and rubrics as new inputs to feed for PLMs. In the rest of this section, we present a two-stage pipeline to score student responses based on the rubrics of the assessments.

## 3.1 Matching Exemplars

Most researchers [29,15,16] reformat their downstream tasks as the masked language modeling task [5], where PLMs are asked to fill the blank [MASK] that is shown in a context with a proper word. For example, by combining a student response and the corresponding rubric, we can construct a new input as "The student responses are that [Response]. The rubric is that [Rubric]. Overall, the response can be grad as [MASK] point.", where each underlined word is filled with the actual input, and the rest parts of this example are called template. Ideally, we expect the PLMs to fill the blank [MASK] with the grade of the

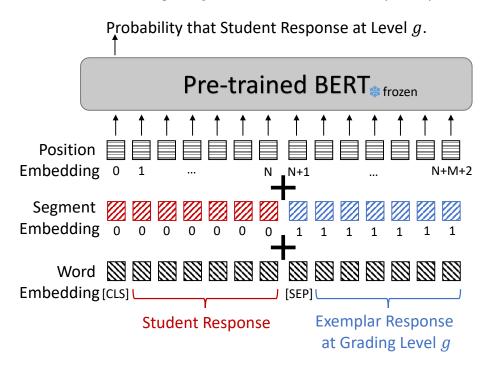


Fig. 1. Next sentence prediction as matching exemplars. response according to the given scoring rubric. However, our piloting data show that PLMs almost randomly fill in the blank without considering the context. We assert that this strategy is failed for two reasons. Firstly, both the responses and rubrics are so long that PLMs pay less attention to the task definition provided by the template. Secondly, PLM has no hints to bridge the gap between the rubrics and the final grading points.

To some extent, scoring student responses according to a given rubric can be considered a matching process if the rubrics can be refined to each grading level. Following this idea, we propose reforming the scoring procedure as the Next Sentence Prediction (NSP) task [5], which requires PLMs to determine whether the given two texts come from the same context. Specifically, assuming that the rubrics on each grading level are available, we first concatenate the student responses with each rubric's overall grading levels independently. After that, PLMs judge whether each response-rubric pair shows the same context by using the outputs of the NSP task. Finally, each score of a student response will be the grading level indicated by the rubrics that best match it. Since the objective of the NSP task has been clear to PLMs, we no more need to worry about designing templates to let the PLMs know what we require them to do. To this end, we address the first concern of using the MLM task: PLMs cannot easily identify what the scoring task is. On the other hand, to fill the gap between the rubrics and the student responses, we propose replacing the rubrics with exemplars as standards to score the responses. This strategy is necessary because the language style of the student responses dramatically differs from that of the assessment

rubrics, which leads PLMs to constantly predict the rubric-response pair as different (Here is a sample case<sup>1</sup>). Figure 1 demonstrates the proposed method.

Formally, given a pre-defined vocabulary set  $\mathcal{V}$ , a student response  $\mathcal{R} = [w^{(1)},...,w^{(n)},...,w^{(N)}]$  consists of N words  $w^{(n)} \in \mathcal{V}$ , and the g-level grade exemplar response  $\mathcal{E}_g = [w_g^{(1)},...,w_g^{(m)},...,w_g^{(M)}]$  has M words  $w^{(m)} \in \mathcal{V}$ . We adopt a PLM  $f_{\rm nsp}: \mathcal{V}^L \to [0,1]$  with a pre-trained NSP head to automatically score the student response as:

$$score(\mathcal{R}) = \arg\max_{g} f_{nsp}([w_{cls}; \mathcal{R}; w_{sep}; \mathcal{E}_g]),$$
 (1)

where  $w_{\rm cls}, w_{\rm sep} \in \mathcal{V}$  are special words,  $[\cdot; \cdot]$  indicates the concatenating operation over the input words, and L = M + N + 2 in this case. According to Equation 1, we achieve automatic scoring under the zero-shot scenarios (without training) by aligning the scoring process into an exemplar matching process and releasing the PLMs' ability learned from the pre-training stage.

#### 3.2 Zero Grade Identifier

However, the above strategy raises a new challenge in obtaining fine-grained exemplars for each grading level. We simultaneously cope with this challenge by decoupling the rubric and the perfect response of an assessment item. Particularly, the perfect rubric can be separated into several points, where each point reflects a grading level. Therefore, it is easy to write down the fine-grained rubrics of each grading level. Similarly, we first develop an optimal student response that fits all grading points of the perfect (full score) rubric. Once the optimal response is given, we remove each part of the perfect response to generate exemplar responses at different levels gradually.

Although this strategy generates high-quality responses for non-zero grading levels, it cannot and is also impossible to enumerate all zero-point exemplars since the reasons that the responses receive high scores are limited (exhaustive). In contrast, the zero-point responses can be various (in-exhaustive). To reduce invalid scoring because of missing zero-point exemplars, we introduce a pre-stage before the above method so that we find out those zero-point responses in an early stage and let Equation 1 focus on how many points the responses reward. Recall that PLMs represent inputs as vectors in their interior, carrying rich semantic information. Since the zero-point response is very different from the perfect response, the distance between the vectors of them should be far away. Thus, PLMs generate vectors to represent both the perfect exemplar response and the given student response and measure the distance between these two vectors with the Cosine similarity. If the cosine similarity between them is smaller

<sup>&</sup>lt;sup>1</sup> Rubric: Student can specify Sam's claim and corresponding evidence, and explain the relationship between the claim and evidence properly. Response: Sam's claim is that gas particles float on the top in the box. The evidence is that bubbles in soda water float to the top. As gas particles and bubbles, all refer to air. Sam can infer that air in the box can also float on the top just as air in soda water does.

than a threshold, we directly grade them with zero point. The threshold is calculated by averaging the cosine similarities of the zero-score and the one-score exemplar response to the perfect exemplar response.

Theoretically, given a PLM  $f_{\rm emb}: \mathcal{V}^L \to \mathbb{R}^d$  that maps a piece of L-length text into a d-dimensional space, we collect the representations of the student response and the g-level exemplar  $\mathbf{z}_{\rm R} = f_{\rm emb}(\mathcal{R})$  and  $\mathbf{z}_g = f_{\rm emb}(\mathcal{E}_g)$ , respectively. We determine whether the response  $\mathcal{R}$  belongs to 0-level grade by:

$$zero\_score(\mathcal{R}) = \begin{cases} Yes, cos(\mathbf{z}_{R}, \mathbf{z}_{g=2}) < \theta, \\ No, otherwise, \end{cases}$$
 (2)

where

$$\theta = \frac{\cos(\mathbf{z}_{g=0}, \mathbf{z}_{g=2}) + \cos(\mathbf{z}_{g=1}, \mathbf{z}_{g=2})}{2},$$
(3)

 $\cos(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}\mathbf{y}^{\mathsf{T}}}{||\mathbf{x}||_2||\mathbf{y}||_2}$  is the cosine distance among vectors, and  $\mathbf{z}_{g=i} \in \mathbb{R}^d$  is the embedding of the *i*-th grading level exemplar, . The combination of Equation 2 and 3 identifies most zero-score responses without enumerating all possible zero-level exemplars. Since Equation 2 and 3 are non-parametric, we can identify zero-score responses under the zero-shot setting.

## 4 Experiment

This section aims to quantitatively justify the following three research questions (**RQ**): (1) How accurate is the proposed method in scoring student responses? (2) Can the performance of the proposed method be further improved if a few sample responses are available? (3) How does the quality of sample responses affect the method performance under the few-shots setting?

## 4.1 Setup

Dataset We choose a subset of an existing dataset of argumentation items [36] for our experiments. The dataset originally consists of eight written response items sharing the same context regarding gases. The items require varied levels of cognitive demands aligned with a learning progression of argumentation in science [22]. Students' responses were scored based on their performance on claim stating, evidence clarifying, and warrants using. According to the different cognitive demands of each item, the scoring rubrics varied between 2 to 4 levels. Meanwhile, each item is designed with different levels of complexity, diversity (refers to the three-dimensional science learning requirement), and structure (refers to the learning progression of scientific argumentation [10]. In this study, we select three items (G4, G5, and G6) sharing the same diversity, structures, and rubrics level but distributing in two complexity levels as our downstream tasks. Overall, the dataset contains 2081 labeled responses (770 for item G4, 669 for item G5, and 642 for item G6) from 931 students of grades 5 to 8.

Dataset Splitting To mimic a few-shot setting, we only leave three samples at each grading level to answer **RQ2** and **RQ3**. Here, no valid set remains because tuning the hyper-parameters for each item is not encouraged. Since almost the entire dataset is considered the test set, the K-folds setting is unnecessary.

Metric We calculate the Cohen's kappa (Kappa) and F1 scores to measure the performance of the machine learning models for auto-scoring. Cohen's kappa is one of the standard rater agreement indices to quantify levels of agreement between computer scoring and human expert scoring [2]. Kappa value ranges from -1 to 1, where 1 indicates a strong agreement between the human score and the machine score and 0 refers to an opposite agreement. Typically, a machine learning model can be accepted if it reaches around 0.4 Kappa score [39]. We also report the F1 score, the weighted harmonic mean of Recall and Precision, to evaluate the model performance [24].

Baseline We compare MeNSP with some baseline methods to measure the effectiveness of MeNSP. Under the zero-shot setting, we perform Random strategy as our baseline, which randomly scores the student responses from 0 to 2 (three grading levels). Under the few-shots setting, we follow the previous studies [36] that performed machine learning in the automatic scoring task. Specifically, we choose some popular ensemble models as the baselines, including Gradient-Gased Decision Tree (GBDT), Random Forest with Decision Tree (RFDT), and the simple Voting strategy (Vote) over five basic models (e.g., Naive Bayes, Decision Tree, Logistic Regression, Multilayers Perceptron, and Support Vector Machine). All baselines make decisions based on the TF-IDF [26] scores of words presented in the student responses.

**Exemplar Design** We first manually develop exemplars for each task according to the rubric of the highest level (level 2) [10]. Human experts are involved to ensure that the level 2 exemplars contain all elements of a perfect argument to the greatest extent. Then, we delete elements level-by-level and adapt sentences to meet the rubrics of level 1 and level 0. These exemplars are the prompts used in the zero-shot experiment.

Sample Design To examine the model performance training with few samples, we use two strategies to generate samples for the few-shots tuning: (1) We randomly select student responses for each level from the dataset. (2) We use ChatGPT<sup>2</sup> to generate new responses based on the exemplar at each grading level and then conduct manual inspection and adjustment to ensure the machine-generated responses meet the rubrics.

#### 4.2 Results

To reduce the uncertainty of experiments, we run through our experiment codes over five random seeds, including data splitting and few-shots tuning, and report both the mean and standard deviation of each metric on each item in Table 1.

<sup>&</sup>lt;sup>2</sup> ChatGPT is available at https://chat.openai.com/chat.

Shot	Sample	Model	G4		G5		G6	
			Kappa (%)	F1 (%)	Kappa (%)	F1 (%)	Kappa (%)	F1 (%)
0	-	Random	$-0.2_{\pm 3.4}$	$32.9_{\pm 1.9}$	$-1.1_{\pm 2.2}$	$35.8_{\pm 1.5}$	$-3.3_{\pm 3.6}$	$30.6_{\pm 2.0}$
		MeNSP	$30.3_{\pm 0.3}$	$54.2_{\pm0.2}$	$57.2_{\pm 1.1}$	$81.1_{\pm 0.5}$	$34.5_{\pm 0.3}$	$57.0_{\pm 0.3}$
1	Random	RFDT	$-1.1_{\pm 2.9}$	$23.8_{\pm 5.5}$	$6.5_{\pm 9.4}$	$50.0_{\pm 14.8}$	$-1.1_{\pm 4.5}$	$24.1_{\pm 11.3}$
		GBDT	$5.6_{\pm 7.2}$	$34.7{\scriptstyle\pm7.3}$	$6.1_{\pm 16.2}$	$43.6{\scriptstyle\pm17.9}$	$3.3_{\pm 9.9}$	$34.6{\scriptstyle\pm11.7}$
		Vote	$3.1_{\pm 11.2}$	$33.6_{\pm9.1}$	$20.9_{\pm 12.2}$	$62.0_{\pm 8.9}$	$7.3_{\pm 9.5}$	$36.5_{\pm 14.1}$
		MeNSP	$35.9_{\pm 3.2}$	$57.9_{\pm 2.2}$		$74.3_{\pm 6.1}$	$25.7_{\pm 4.3}$	$53.2_{\pm 3.6}$
	Manual	RFDT	$0.1_{\pm 0.1}$	$26.3_{\pm 0.1}$	$0.3_{\pm 0.3}$	$51.5_{\pm 0.1}$	$0.0_{\pm 0.0}$	$34.2_{\pm 0.0}$
		GBDT	$0.0_{\pm 0.0}$	$26.2{\scriptstyle\pm0.0}$	$1.5_{\pm 1.5}$	$52.5_{\pm 1.0}$	$0.7_{\pm 0.3}$	$34.8_{\pm 0.3}$
		Vote	$7.0_{\pm 1.7}$	$34.0{\scriptstyle\pm1.4}$	$5.6_{\pm 0.7}$	$52.2{\scriptstyle\pm0.4}$	$-0.14_{\pm 0.0}$	$34.2{\scriptstyle\pm0.0}$
		MeNSP	$35.5_{\pm 7.1}$	$59.1_{\pm1.8}$	$27.0_{\pm 4.5}$	$56.5_{\pm 0.5}$	$30.9_{\pm 1.8}$	$56.8_{\pm 0.6}$
3	Random	RFDT	$5.5_{\pm 6.2}$	$33.3_{\pm 5.9}$	$13.2_{\pm 14.4}$	$49.8_{\pm 19.8}$	$3.6_{\pm 2.2}$	$33.1_{\pm 11.2}$
		GBDT	$4.0_{\pm 7.2}$	$33.2{\scriptstyle\pm7.4}$	$14.0_{\pm 16.0}$	$49.6{\scriptstyle\pm19.9}$	$13.6_{\pm 3.6}$	$46.2{\scriptstyle\pm1.7}$
		Vote	$15.6_{\pm 3.4}$	$44.8{\scriptstyle\pm2.0}$	$20.9_{\pm 16.7}$	$57.4{\scriptstyle\pm10.8}$	$13.1_{\pm 4.9}$	$44.1{\scriptstyle\pm6.9}$
		MeNSP	$37.5_{\pm 3.3}$	$59.0_{\pm2.0}$		$78.7_{\pm 3.6}$	$27.8_{\pm 6.7}$	$53.4_{\pm 4.6}$
	Manual	RFDT	8.1 <sub>±3.9</sub>	$34.2_{\pm 3.7}$	$0.7_{\pm 0.4}$	$51.6_{\pm 0.1}$	$0.0_{\pm 0.0}$	$34.2_{\pm 0.0}$
		GBDT	$17.1_{\pm 0.4}$	$45.3{\scriptstyle\pm0.3}$	$7.6_{\pm 0.2}$	$53.2_{\pm 0.1}$	$1.03_{\pm 0.0}$	$35.64_{\pm 0.0}$
		Vote	$14.8_{\pm 1.0}$	$43.9{\scriptstyle\pm1.0}$	$8.1_{\pm 2.1}$	$52.5{\scriptstyle\pm1.1}$	$7.83_{\pm 0.6}$	$40.07{\scriptstyle\pm0.4}$
		MeNSP	$33.8_{\pm 1.2}$	$56.2_{\pm 0.9}$	$24.8_{\pm 1.1}$	$58.2_{\pm 1.4}$	$30.1_{\pm 0.5}$	$56.9_{\pm 0.3}$

**Table 1.** Machine automatically scoring performance.

RQ1: MeNSP effectively scores responses based on context and the exemplars without training (zero-shot). Under the zero-shot setting, all the Kappa values of the three items increase from negative to more than 0.30, and the F1 scores also increase among the three items, which indicates the effectiveness of MeNSP in the zero-shot condition (compared with Random baseline). One of the items (G5) has its Kappa exceed 0.50 with MeNSP, which reaches an acceptable benchmark of the trained model. Although MeNSP surpluses a 0.4 Kappa score only on G5, it still achieves at least a 0.3 Kappa score on the other two items, showing the great potential of MeNSP.

RQ2: MeNSP can be improved by training on a few sample responses on easier items (few-shots). To discuss with the few-shots setting, MeNSP performs better than itself on item G4, with the Kappa of random MeNSP increasing from 0.30 (zero-shot) to 0.36 (1-shot), and to 0.38 (3-shots). However, increasing the number of sample responses does not improve MeNSP further on G5 and G6. Although MeNSP's performance fails on two items, it still has a higher average Kappa and lower standard deviation than all baselines under the 1-shot and 3-shots settings on all three items. This indicates that MeNSP performs better and is more stable than other models.

Considering the lower complexity level of G4 (level 2) than G5 and G6 (level 3), we first conclude that MeNSP's performance can be improved by training on some responses on the low complexity items. However, for items with high complexity, we argue that the performance of MeNSP heavily relies on the characteristics of the sample responses provided for training. According to our ex-

periments, we observe an improvement of Kappa on G5 from 0.59 to 0.61 with random seed 55301 and from 0.57 to 0.60 with random seed 9, and a gain of Kappa on G6 from 0.34 to 0.35 with random seed 45983. A possible speculation of the successful improvement with these random seeds is that the sample responses selected for training are diverse enough<sup>3</sup> to help MeNSP to learn complete and unbiased knowledge about the general scoring rubrics. This speculation also aligns with the finding of RQ3.

RQ3: The quality of the provided samples affects the few-shots learning performance. We compare MeNSP on item G4 with different sample-gathering strategies (random or manual). The results show that with the 1-shot setting, both strategies lead MeNSP to a similar performance, with the average Kappa close to 0.36. However, with the 3-shot setting, the random strategy performs better than the manual strategies, with the average Kappa increase from 0.34 (manual) to 0.38 (random). As the random strategy extracts samples from the real student responses, the manual strategy provides augmented student responses similar to the exemplars<sup>4</sup>. Therefore, combined with the observation of improvements with three random seeds on G5 and G6 discussed earlier, we summarize that the diversity of the given samples is a potential factor that impacts the efficiency of few-shots learning.

#### 5 Conclusion and Discussion

In this study, we develop a zero-shot approach (MeNSP) to score student responses automatically. We propose three research questions and examine MeNSP's performance with 1-shot and 3-shot scenarios. Through experiments, we prove the effectiveness of MeNSP on automatically scoring student responses based on exemplars without training. We also find that increasing training samples can improve MeNSP's performance on items with lower complexity. However, the quality of the improvement is related to the characteristics of the training samples, for example, the diversity of the sample responses. Given that our goal of this study is to preliminary demonstrate the feasibility of MeNSP and the prompt learning methods in automatic scoring science assessment, the machine scoring accuracy may be used for low-stake formative assessment practices. Future research should improve the performance of MeNSP and use varying datasets to verify the approach.

<sup>&</sup>lt;sup>3</sup> G5 random samples at grade 2: (1) if the gas bubbles in soda go to the top, then so should the gas particles in the box. (2) Sam could use the Gas bubble idea to support his claim. Gas bubbles are made up of particles. The gas bubbles float to the top of the can. So, the particles in the box will rise to the top also. (3) The gas particles go to the top, like in soda.

<sup>&</sup>lt;sup>4</sup> G4 manual samples at grade 1: (1) The evidence suggests that the air in the balloon has spread throughout. (2) Charlie's claim is right because the air in the balloon spreads everywhere. (3) That air in the balloon spreads everywhere supports claim.

# Acknowledgement

The study was funded by National Science Foundation(NSF) (Award # 2101104, 2138854, PI: Zhai). Any opinions, findings, conclusions, or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the NSF.

## References

- 1. Amerman, H., et al.: Does transformer deep learning yield more accurate sores on student written explanations than traditional machine learning? In: AERA
- 2. Bejar, I.I.: A methodology for scoring open-ended architectural design problems. Journal of Applied Psychology (1991)
- 3. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. Advances in neural information processing systems (2020)
- 4. Council, N.R., et al.: A framework for K-12 science education: Practices, crosscutting concepts, and core ideas. National Academies Press (2012)
- Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: ACL (2019)
- Gao, T., Fisch, A., Chen, D.: Making pre-trained language models better few-shot learners. arXiv preprint arXiv:2012.15723 (2020)
- 7. Gerard, L., Kidron, A., Linn, M.C.: Guiding collaborative revision of science explanations. Int. Journal of Computer-Supported Collaborative Learning (2019)
- 8. Harris, C.J., et al.: Designing knowledge-in-use assessments to promote deeper learning. Educational measurement: issues and practice (2019)
- Haudek, K.C., et al.: What are they thinking? automated analysis of student writing about acid-base chemistry in introductory biology. Life Sciences Education (2012)
- 10. Haudek, K.C., Zhai, X.: Exploring the effect of assessment construct complexity on machine learning scoring of argumentation (2021)
- 11. Lee, H.S., et al.: Automated text scoring and real-time adjustable feedback: Supporting revision of scientific arguments involving uncertainty. Science Education (2019)
- 12. Litman, D.: Natural language processing for enhancing teaching and learning. In: Thirtieth AAAI conference on artificial intelligence (2016)
- 13. Liu, O.L., et al.: Automated scoring of constructed-response science items: Prospects and obstacles. Educational Measurement: Issues and Practice (2014)
- 14. Liu, P., et al.: Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. arXiv preprint arXiv:2107.13586 (2021)
- 15. Liu, X., et al.: Gpt understands, too. arXiv preprint arXiv:2103.10385 (2021)
- 16. Liu, X., et al.: P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. arXiv preprint arXiv:2110.07602 (2021)
- 17. Lu, Y., et al.: Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. arXiv preprint arXiv:2104.08786 (2021)
- 18. Maestrales, S.e.a.: Using machine learning to score multi-dimensional assessments of chemistry and physics. Journal of Science Education and Technology (2021)
- Mayer, C.W., Ludwig, S., Brandt, S.: Prompt text classifications with transformer models! an exemplary introduction to prompt-based learning with large language models. Journal of Research on Technology in Education (2022)

- Nehm, R.H., Ha, M., Mayfield, E.: Transforming biology assessment with machine learning: automated scoring of written evolutionary explanations. Journal of Science Education and Technology (2012)
- 21. Omizo, R., Meeks, M., Hart-Davidson, W.: Detecting high-quality comments in written feedback with a zero shot classifier. In: ACM ICDC (2021)
- 22. Osborne, J.F., et al.: The development and validation of a learning progression for argumentation in science. Journal of research in science teaching (2016)
- Pellegrino, J.W.: Proficiency in science: Assessment challenges and opportunities. Science (2013)
- Powers, D.M.: What the f-measure doesn't measure: Features, flaws, fallacies and fixes. arXiv preprint arXiv:1503.06410 (2015)
- 25. Riordan, B.e.a.: An empirical investigation of neural methods for content scoring of science explanations. In: Proceedings of the fifteenth workshop on innovative use of NLP for building educational applications (2020)
- Salton, G., Buckley, C.: Term-weighting approaches in automatic text retrieval. Information processing & management (1988)
- 27. Schick, T., Schütze, H.: It's not just size that matters: Small language models are also few-shot learners. arXiv preprint arXiv:2009.07118 (2020)
- 28. Schick, T., Schütze, H.: Exploiting cloze questions for few-shot text classification and natural language inference. Computing Research Repository (2020)
- 29. Shin, T., et al.: Autoprompt: Eliciting knowledge from language models with automatically generated prompts. arXiv preprint arXiv:2010.15980 (2020)
- Su, Y., et al.: On transferability of prompt tuning for natural language processing.
   In: NACL. pp. 3949–3969 (2022)
- 31. Uhl, J.D., et al.: Introductory biology undergraduate students' mixed ideas about genetic information flow. Biochemistry and Molecular Biology Education (2021)
- 32. Vu, T., et al.: Spot: Better frozen model adaptation through soft prompt transfer. arXiv preprint arXiv:2110.07904 (2021)
- 33. Wolfe, E.W., Wendler, C.L.W.: Why should we care about human raters? Applied Measurement in Education (2020)
- 34. Wu, X., et al.: A survey of graph prompting methods: Techniques, applications, and challenges. arXiv preprint arXiv:2303.07275 (2023)
- 35. Zhai, X.: Practices and theories: How can machine learning assist in innovative assessment practices in science education. Journal of Science Education and Technology (2021)
- 36. Zhai, X., Haudek, K.C., Ma, W.: Assessing argumentation using machine learning and cognitive diagnostic modeling. Research in Science Education (2022)
- 37. Zhai, X., Krajcik, J., Pellegrino, J.W.: On the validity of machine learning-based next generation science assessments: A validity inferential network. Journal of Science Education and Technology (2021)
- 38. Zhai, X., Shi, L., Nehm, R.H.: A meta-analysis of machine learning-based science assessments: factors impacting machine-human score agreements. Journal of Science Education and Technology (2021)
- 39. Zhai, X., Yin, Y., Pellegrino, J.W., Haudek, K.C., Shi, L.: Applying machine learning in science assessment: a systematic review. Studies in Science Education (2020)
- 40. Zhang, M., et al.: Automatic short math answer grading via in-context meta-learning. arXiv preprint arXiv:2205.15219 (2022)
- 41. Zhong, R., Lee, K., Zhang, Z., Klein, D.: Adapting language models for zero-shot learning by meta-tuning on dataset and prompt collections. In: EMNLP (2021)