Hanling Wang

Xi'an Jiaotong-Liverpool University School of Advanced Technology Suzhou, Jiangsu, China Hanling.Wang21@student.xjtlu.edu.cn

Kexin Chen Xi'an Jiaotong-Liverpool University School of Advanced Technology Suzhou, Jiangsu, China kexin.chen21@student.xjtlu.edu.cn

Yiwei Li Xi'an Jiaotong-Liverpool University School of Advanced Technology Suzhou, Jiangsu, China Yiwei.Li21@student.xjtlu.edu.cn

Abstract

With the advancement of educational technology, automatic assessment systems are becoming increasingly essential, particularly for grading short-answer questions. However, due to the inherent ambiguity and complexity of language, automatic grading of shortanswer questions remains a challenge. Traditional grading methods are often time-consuming and subjective, highlighting the need for efficient, objective, and feedback-driven solutions. This paper proposes an innovative approach to automatic short answer grading (ASAG) utilizing large language models (LLMs). We introduce a specialized design for crafting questions and corresponding answers named Key Point Scoring Framework (KPSF) which significantly enhances the model's performance in ASAG tasks and improves the flexibility and objectivity of assessments. Moreover, we incorporate Prompt Dynamic Adjustment (PDA) that continuously refines the grading process, effectively handling ambiguous student responses while ensuring reliable results. To evaluate our approach, we develop a multidisciplinary dataset and incorporate real-world dataset from actual exams. The experimental results demonstrate that our ASAG approach provides educators with a highly efficient, flexible and accurate tool for short-answer assessments, indicating a significant advancement in automatic grading technology.

Banghao Chi

University of Illinois Urbana-Champaign College of Liberal Arts & Sciences Urbana, IL, USA banghao2@illinois.edu

Di Wu

Xi'an Jiaotong-Liverpool University School of Advanced Technology Suzhou, Jiangsu, China Di.Wu23@student.xjtlu.edu.cn

Hanyan Niu Xi'an Jiaotong-Liverpool University School of Advanced Technology Suzhou, Jiangsu, China Hanyan.Niu23@student.xjtlu.edu.cn

CCS Concepts

• Computing methodologies \rightarrow Natural language processing; • Applied computing \rightarrow Education; Computer-assisted instruction;

Keywords

Large Language Models, Key Point Scoring, Natural Language Processing, Prompt Engineering, Automatic Grading

ACM Reference Format:

Hanling Wang, Banghao Chi, Yufei Wu, Kexin Chen, Di Wu, Songning Liu, Yiwei Li, Hanyan Niu, and Xiaohui Zhu. 2025. LLMarking: Adaptive Automatic Short-Answer Grading Using Large Language Models. In Proceedings of the Twelfth ACM Conference on Learning @ Scale (L@S '25), July 21–23, 2025, Palermo, Italy. ACM, New York, NY, USA, 11 pages. https: //doi.org/10.1145/3698205.3729551

1 Introduction

In recent years, recent innovations in educational technology have transformed many aspects of teaching and assessment. One area that garners significant attention is automatic grading, particularly in short-answer assessments. Automatic Short Answer Grading (ASAG) has emerged as a powerful tool that utilizes computer algorithms to analyze and evaluate student responses to open-ended questions, offering substantial benefits in both efficiency and quality of assessment. First, ASAG improves grading efficiency, particularly in large classrooms with high teacher workloads, by automating the process and allowing teachers to focus on more impactful activities [21, 29]. This shift increases productivity and enhances the educational experience [19]. Second, ASAG ensures greater consistency in grading, reducing variability from human bias, fatigue, or interpretation differences [11, 28]. Automatic systems provide equitable assessments for all students, minimizing these impacts

Yufei Wu

Xi'an Jiaotong-Liverpool University School of Advanced Technology Suzhou, Jiangsu, China Yufei.Wu21@student.xjtlu.edu.cn

Songning Liu

Xi'an Jiaotong-Liverpool University School of Advanced Technology Suzhou, Jiangsu, China Songning.Liu21@student.xjtlu.edu.cn

Xiaohui Zhu* xiaohui.zhu@xjtlu.edu.cn Xi'an Jiaotong-Liverpool University School of Advanced Technology Suzhou, Jiangsu, China

^{*}Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org. La@S '25. Palerno. Italy

L@S '25, Palermo, Italy

^{© 2025} Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 979-8-4007-1291-3/2025/07 https://doi.org/10.1145/3698205.3729551

[27]. Finally, integrating large language models (LLMs) into ASAG adds explainability, offering clear justifications for scoring decisions and helping students understand their mistakes and improve [14]. This transparency fosters trust in the grading process and supports a learning-oriented approach.

Despite its advantages, ASAG presents significant challenges due to the complexity of natural language understanding, which requires a nuanced analysis of student responses [28]. One major limitation is that many frameworks provide only a general score without offering detailed feedback on how specific points in the reference answer are addressed. This lack of granular feedback makes it difficult for students to understand their errors, reducing learning opportunities and hindering meaningful progress. Additionally, ASAG struggles to accurately match student responses with reference answers due to variations such as synonyms, paraphrasing, and implied meanings [7]. While human graders can recognize these nuances, current ASAG systems often fail to capture them, leading to inaccurate assessments. Another challenge is the inability of ASAG models to handle incomplete or ambiguous responses, such as those with missing information or indirect reasoning, resulting in inconsistent grading [25]. Furthermore, current studies assess a limited range of LLMs, restricting the generalizability and effectiveness of ASAG systems. Broader evaluation across diverse LLMs is necessary to address these weaknesses and improve feedback mechanisms.

To address these challenges, this paper proposes LLMarking algorithm, which aims to improve both the scoring accuracy and consistency of automated grading. The key objectives of our work are fourfold: (1) enhance the grading system's feedback mechanism to offer detailed feedback on the score judgment, allowing students to better understand the rationale behind their grades; (2) develop a more flexible and fair grading method that can adapt to various types of open-ended questions and more accurately match student responses with reference answers; (3) introduce mechanisms to handle incomplete or ambiguous student responses more effectively; and (4) give comprehensive and broader performance evaluation of LLMs, ensuring that a wider range of LLMs are assessed to uncover strengths and weaknesses in their grading capabilities. By addressing these key areas, this research aims to advance the effectiveness of ASAG systems, making them more reliable and useful in educational contexts.

Our work introduces a Key Point Scoring Framework (KPSF), in which the reference answer is manually divided into key points to create a detailed rubric for evaluating student responses. This structure allows LLMs to assess student answers based on clear and well-defined criteria, making the evaluation process more transparent and aligned with the expectations of the reference answer. In addition, LLMarking algorithm incorporates Prompt Dynamic Adjustment (PDA) mechanism which enhances the system's adaptability to ambiguous responses. This dynamic adjustment allows the algorithm to account for variations in phrasing, paraphrasing, and other linguistic nuances that might otherwise lead to inaccurate assessments. Moreover, we evaluate different types of LLMs to identify the best-performing model for auto-grading, ensuring that the best LLM is selected in the research to enhance reliability. Integrating these features ensures that students receive detailed, meaningful feedback, promoting a more fair and equitable grading

process and learning improvement by offering insights into how they can enhance their responses.

The main contributions of our study are as follows:

- Utilization of Leading LLMs: Our work leverages the latest advancements in LLMs to improve ASAG. By evaluating and employing leading LLMs in combination with the KPSF, we enhance the model's ability to provide detailed feedback. This approach allows students to clearly understand how their answers align with the reference answer and where they go wrong, with specific reasons for each point awarded or lost, thereby promoting a deeper understanding of their mistakes and guiding improvement.
- Key Point Scoring Framework (KPSF): We develop a KPSF for ASAG that breaks down the reference answer into key points, allowing accurate and flexible matching of diverse student responses. This improves model performance and addresses the limitations of traditional scoring methods.
- Prompt Dynamic Adjustment (PDA): PDA continuously refines the grading process, effectively managing ambiguous student responses and ensuring reliable results across various subjects. This mechanism addresses the weakness of LLMs in dealing with incomplete or ambiguous data, enhancing the accuracy and fairness of scoring.
- Dataset and Comprehensive Evaluation:

We assess our framework by developing a multidisciplinary dataset based on real exam papers, demonstrating its effectiveness with leading LLMs. Additionally, we plan to make our dataset publicly available, providing a valuable open resource for further research in ASAG. This broad evaluation will help to enhance the grading capabilities of ASAG systems. Code and data are available at https://github.com/2024-Surf-LLMarking/LLMarking.

2 Related Work

2.1 Automatic Short Answer Grading

ASAG is an important advancement in educational technology, designed to reduce the effort required for manual grading while maintaining consistency and objectivity in assessments. The early development of ASAG relies on rule-based systems and keywordmatching methods. For example, Burrows et al. [3] outlined initial trends in ASAG, where concept mapping was used to compare responses by measuring similarity to reference answers. However, these approaches are limited by their inability to capture semantic meaning and rely heavily on surface-level features, leading to inaccuracies in grading. To address the shortcomings of rule-based approaches, information retrieval techniques are introduced. Pulman [24] applied domain-specific patterns to extract key details, but the lack of flexibility remains a challenge. Corpus-based methods later incorporate statistical and semantic features, improving adaptability but still failing to fully capture the nuances of student responses [20].

Machine learning models, particularly support vector machines (SVMs) and bag-of-words approaches [12], offer greater flexibility in ASAG. These models are widely used due to their effectiveness in various scenarios. However, their reliance on fundamental features such as n-grams restricts their ability to fully capture contextual meaning in student responses. While these methods perform well in many cases, they often struggle when faced with complex or nuanced answers that require a deeper understanding of context and semantics [9].

Deep learning models have significantly enhanced ASAG by providing a deeper contextual understanding, improving the ability to grasp word dependencies and effectively handling complex linguistic structures. These advancements have made automated grading more accurate and reliable. De Mulder et al. [6] demonstrated the power of recurrent neural networks (RNNs) in processing sequential data, while Cheng et al. [5] extended this capability using Long Short-Term Memory networks (LSTMs), which are particularly effective at capturing long-range dependencies in text. Further refining these approaches, Mueller and Thyagarajan [22] introduced a Siamese LSTM architecture for paired sequence comparison, enhancing the ability to assess semantic similarity in student responses. Building on this, Kumar et al. [16] proposed a Siamese biLSTM combined with a Sinkhorn distance pooling layer to further improve sequence comparison and grading accuracy. Despite these notable advancements, deep learning models still face challenges in ensuring grading consistency, particularly when dealing with ambiguous, misleading, or highly varied student answers, which can introduce subjectivity and affect overall reliability.

Large Language Models (LLMs), such as Lamma-3 and GPT-4, have significantly improved upon previous deep learning approaches by better capturing nuanced relationships in language. These models effectively address the ambiguities that traditional RNNs and LSTMs struggled with, leading to fewer grading inconsistencies and high accuracy in ASAG. Their advanced natural language processing capabilities enable them to understand context more deeply, making them particularly useful for complex and varied student responses. Yoon [32] demonstrated this by employing one-shot prompting with GPT-3.5 to extract key phrases from student answers, showcasing its effectiveness in identifying relevant information. Building on this, Hackl et al. [10] further highlighted GPT-4's consistency in text evaluation, showing a high intraclass correlation that indicates reliable grading performance. Despite these improvements, LLMs continue to struggle with grading consistency, particularly when encountering highly ambiguous or misleading responses, as highlighted by Chang et al. [4].

2.2 Dynamic Prompting

Dynamic Prompting has emerged as a pivotal technique that enhances model performance and adaptability by tailoring prompts in real-time to align with specific tasks or user interactions. This approach transcends the limitations of static prompts, enabling models to respond more effectively to diverse inputs and contexts. Prior research has extensively explored the application of Dynamic Prompting across various domains, demonstrating its versatility and efficacy. Yang et al. [31] introduced a unified dynamic prompt tuning strategy, which dynamically determines multiple factors such as prompt position, length, and representation based on specific tasks and instances. Their work demonstrates that optimizing prompt placement can capture additional semantic information, which traditional prefix or postfix prompt tuning methods fail to encapsulate. However, the disadvantage of it is that this

method mainly focuses on the optimization of prompt structure without fully addressing context-specific details in real-time responses. Zhao et al. [33] proposed a Dynamic Prompt Adjustment framework to address knowledge forgetting in multi-label classincremental learning. Their approach integrates an improved data replay mechanism alongside prompt loss regularization, enabling adaptive prompt modification for evolving learning environments. Yet, the main limitation of this method is its complexity, as it requires substantial computational resources to implement the data replay mechanism effectively, which may not be feasible in all practical applications. Additionally, Kamesh [26] introduced Adaptive Prompting, a framework designed to enhance reasoning capabilities in LLMs through real-time adjustments to prompt structures and validation mechanisms. However, this research lacks a comprehensive strategy for dealing with vague or misleading student responses, which can lead to inconsistencies in grading results.

Building on these prior studies, our proposed algorithm, *LL-Marking*, enhances existing LLM-based methods by introducing structured scoring mechanisms. In particular, we segment the reference answer into labeled components, allowing LLMs to assess student responses using clear criteria. A PDA mechanism, inspired by Fleiss' Kappa, flags ambiguous answers and adjusts prompt dynamically, improving reliability. This blend of structured scoring and prompt dynamic adjustment boosts grading accuracy, making *LLMarking* an effective tool for ASAG.

3 Method

In this section, we present the *LLMarking* algorithm to grade short answers. We introduce the Key Point Scoring Framework (KPSF), which breaks reference answers into labeled points for consistent grading. The section also discusses the datasets used, including cross-subject and real-world exam datasets. We then explain how LLMs assess student responses and describe prompt dynamic adjustment (PDA), which refines prompts to ensure grading accuracy through iterative feedback.

3.1 Key Point Scoring Framework

The grading of short-answer questions requires precise criteria to ensure objectivity and consistency. We design a point-based system in which each important aspect of the answer is assigned a specific score. By breaking down the reference answer into specific labeled scoring points, the LLMs can evaluate students' answers against clear, predefined criteria. Also, this structured format ensures that all responses are evaluated against the same standards, maintaining consistency across various students' answers. Separated points create a clear record of what is assessed and why a particular score is given, providing transparency in the grading process. We develop a label-based reference answer format where points are assigned to specific aspects:

<Point:Mark> specific aspect of answer <Point:Mark>

Each <Point>represents a distinct criterion or detail required in the response, with 'Mark' indicating the score assigned for each correctly addressed point. The use of labels creates a clear separation between key points, making it easier for the model to identify and





evaluate them. KPSF is manually constructed by human annotators, while the evaluation process based on these key points is automated using LLMs. This design ensures both accuracy in rubric definition and scalability in grading.

3.2 Data Collection

In our experiments, we use two types of datasets: a cross-subject question dataset and a real-world exam dataset. Each dataset consists of four components: the question, reference answer, student response, and instructor-assigned score. The cross-subject dataset evaluates the model's ability to generalize across different subjects, while the real-world exam dataset tests whether *LLMarking* can adapt to practical grading scenarios.

Cross-subject Dataset: This dataset comprises independently selected questions from Computer Science (CS), Artificial Intelligence (AI), and Finance (FIN), all curated by subject-matter experts. The questions, sourced from textbooks, online resources, and academic publications, are designed to be unambiguous and suitable for automated evaluation. Each subject includes 8 standalone questions, with the number of evaluation points varying across subjects: 64 points for CS, 21 for AI, and 28 for FIN and every question was answered by 10 different students.

To ensure consistent grading, the reference answers with labeled key points and assigned marks are manually prepared by human annotators, and student responses are independently scored by two instructors based on a standardized rubric. For example, a question about the Software Development Life Cycle may have the following reference answer:

What are the key phases of the Software Development Life Cycle (SDLC)?

<Point1:2>Requirement Gathering <Point2:2>Collecting requirements from stakeholders

<Point3:2>System Analysis and Design

This dataset is used to test the general performance of LLMs in evaluating cross-subject answers based on a KPSF. It helps assess the model's ability to understand and grade responses across different subjects, ensuring consistent and accurate evaluation. The structured reference answers allow for testing the model's ability to handle diverse topics and answer formats.

Real-World Exam Dataset: To evaluate *LLMarking* under practical conditions, we collect a dataset from a real-world computer science exam that includes 10 text-based questions and responses from 40 students. This is a complete exam administered under standard test conditions, with students are required to complete all questions within a fixed time limit.

Unlike the cross-subject dataset, where grading follows a strict predefined rubric with multiple reviewers, this dataset reflects that real-world grading practices involve more flexibility and subjective judgment. Here, a single instructor assigns scores based on the official marking criteria. The reference answers still follow a keypoint format but incorporate variations(cases) commonly observed in actual student responses. For example, a question about project management and risk mitigation can be graded as follows:

If, for any reason, the project team decided to revisit the Software
Specification stage, what measures should be in place to reduce
the negative impacts on the project?
<point1_case1:1>Limit the revisit duration</point1_case1:1>
<point1_case2:1>Reduce duration</point1_case2:1>
<point1_case3:1>Reduce time of revisit</point1_case3:1>
<point2_case1:1>Reduce the number of revisits</point2_case1:1>
<point3_case3:1>More focused areas on spec</point3_case3:1>

These officially graded student responses and instructor-assigned scores are preserved. This dataset captures the complexities and nuances of real-world grading practices to test the effectiveness and flexibility of *LLMarking* in handling subjective judgment.

3.3 LLMarking Workflow

As shown in Figure 1, the judgment process using LLMs involves a systematic approach to evaluate a student's answer against a reference answer based on predefined criteria. In this workflow, the student's answer is compared to the reference answer, which is broken down into key points. Each key point is assessed with binary feedback, and the PDA refines the process for ambiguous answers, ensuring more accurate and consistent grading. The workflow is detailed as follows:

- **Input Preparation:** The process begins with inputting the **question**, the **corresponding reference answer**, and the **student's response**.
- Answer Extraction: The reference answer is automatically decomposed into key points by model, each representing essential aspects of the concept or question. These key points serve as benchmarks for evaluating the completeness and accuracy of the student's response.
- Point-by-Point Judgement and Feedback: The student's answer is analyzed to determine if it addresses each key point from the reference answer. LLMs perform a detailed comparison, checking for relevant terms, concepts, and explanations that align with the expected answers. For each key point, the model provides a binary judgment—'True' or 'False'—and generates feedback explaining the correctness or deficiencies of the student response. This feedback aims to offer constructive insights for improvement.
- **Prompt Dynamic Adjustment (PDA):** To enhance grading accuracy and consistency, dynamic adjustment is applied to ambiguous answers where the model exhibits low confidence. PDA refines the feedback generation based on the model's confidence level. If the model is confident, it generates feedback directly. If confidence is low, the answer is flagged for manual review. Once the manual judgment is provided, the feedback is updated, and the prompt is dynamically adjusted to enhance future performance. Detailed information is provided in the PDA section.

3.4 Prompt Design

The prompt for auto-grading is carefully crafted to instruct LLMs on how to evaluate a student's answer against a provided question and a reference answer. The main objective is to assess the alignment of the student's response with the reference answer using predefined grading criteria. The prompt structure consists of two parts: **Static Predefined Prompt** and **Dynamic Prompt**, as shown in Figure 2.

3.4.1 **Static Predefined Prompt**. The Static Predefined Prompt includes basic grading rules and standard examples.

Instruction Prompt: This prompt provides a comprehensive overview of the grading process, guiding LLMs through the evaluation of a student's answer. It includes the following components:

• **Basic Instructions**: It outlines the general guidelines and specifies the key elements to consider: the question posed to the student, the reference answer (with key points and marking standards), and the student's actual response.



Figure 2: Prompt Design

- Grading Criteria: Embedded within the reference answer, key points are marked with specific tags (<Point >) to evaluate the student's response. The model compares the student's answer to the reference answer, checking for alignment with the key points and assigning a 'True' or 'False' judgment.
- Feedback Generation: After evaluating each point, the model generates feedback, providing explanations for correct or incorrect answers. This helps students understand their mistakes and areas for improvement.
- Anti-misdirection Requests: A keyword filter blocks manipulative language that could influence grading, ensuring evaluations remain focused on academic merit. For example, statements appealing to emotions or asking for high marks without justifying academic content are prevented. This safeguard ensures the model assesses responses based solely on their correctness and relevance.

General Shot: To improve grading accuracy, the prompt adopts a few-shot learning approach by presenting one or more illustrative examples—each comprising a Question, Reference Answer, Student Answer, and model-formatted Feedback. These examples clarify the grading criteria and expected response structure, enabling the model to develop a more precise understanding of the evaluation process and produce more consistent and reliable assessments. *3.4.2 Dynamic Prompt*. Dynamic Prompt offers flexibility for different scenarios, ensuring that grading can be appropriately adjusted based on the student's answer.

Adapted Shot: This type of shot is not present at the outset of the grading process. Instead, it is dynamically introduced by the PDA mechanism as the grading evolves. The Adapted Shot addresses cases where the model initially struggles to assess difficult or ambiguous student answers. PDA intervenes by generating specific example-based prompts that serve as additional guidance for the model, helping it refines its judgment for similar responses in the future, and improving its performance over time. The content of the Adapted Shot follows the same format as the General Shot.

3.5 Prompt Dynamic Adjustment

Algorithm 1 Dynamic Prompt Adjustment based on Confidence Score

Input: Feedback point context (input_ids), model's output logits Output: Updated dynamic prompt

```
1: logits \leftarrow model(input_ids)
```

- 2: probabilities \leftarrow softmax(logits)
- 3: $P_true \leftarrow probabilities[label_map["True"]]$
- 4: $P_{false} \leftarrow probabilities[label_map["False"]]$
- 5: confidence $\leftarrow \max(P_true, P_false)$ //Model is confident, output feedback
- 6: **if** confidence >confidence_threshold **then**
- 7: Judgement \leftarrow P_true >P_false
- 8: feedback \leftarrow generate_feedback(Judgement)
- 9: **Return:** feedback //Model is not confident, manual judgement needed
- 10: **else**
- 11: manual_judgement ← request_manual_review(input_ids)
- 12: updated_feedback ← generate_feedback(manual_judgement)
- 13: $dynamic_prompt \leftarrow update_prompt(updated_feedback)$
- 14: **Return:** updated_feedback
- 15: end if

To enhance the accuracy and stability of our grading model for objective questions, we implement PDA, as shown in the Algorithm 1. This approach allows for iterative refinement of the model's prompts based on its confidence in the generated feedback.

Specifically, the process starts by evaluating the model's confidence through the output probabilities. If the model exhibits a high level of confidence, the feedback is directly generated based on the prediction. However, when the model's confidence falls below a predefined threshold, the judgment is flagged for manual review. The feedback is updated once the manual judgment is provided, and the prompt is adjusted dynamically to incorporate the new information.

The confidence score is computed by evaluating the output logits, which are transformed into probabilities using the softmax function. Logits represent the raw, unnormalized output scores of the model, which are then normalized to form probabilities. The model calculates the maximum probability between the "True" and "False" labels. If the confidence exceeds the predefined threshold of 0.7, as suggested by [18], the feedback is generated automatically. Otherwise, the model requests manual input for judgment, which is then used to update the feedback and dynamically adjust the prompt.

The model's iterative improvement occurs by incorporating examples where manual judgments were needed. These examples are stored and used in future prompts, allowing the model to learn from the instances of uncertainty and gradually improve its ability to make reliable judgments.

In practice, when the model's confidence is low, the user is shown the student's response and the corresponding reference points, and simply asked to judge it as "True" or "False." Optionally, a short explanation can be added. This lightweight interaction enables fast feedback collection and drives dynamic prompt updates with minimal human effort.

4 Experimental Setup

4.1 Dataset

In our experiments, we utilize two types of datasets: the crosssubject dataset and the real-world exam dataset. The cross-subject dataset assesses model's performance in various disciplines, while the real-world exam dataset examines LLMarking's ability to adapt to practical and real-world conditions. From each dataset, we randomly select 3 questions to serve as examples ("shots"). These selected questions are used to serve as templates for the output of the model. The remaining questions in the dataset are reserved for testing purposes.

4.2 Hardware and Software

We deploy LLM using NVIDIA A100 GPUs (40GB VRAM) for highperformance real-time processing. The setup includes a multicore CPU (32GB RAM), 1TB SSD for fast data handling, and runs on Ubuntu for stability. Python 3.8+, PyTorch 2.3, and CUDA 12.1 provide GPU acceleration. FastAPI manages API calls efficiently, while vLLM supports asynchronous inference to improve throughput [17]. Model handling is facilitated by Transformers and Modelscope for optimal integration of pre-trained models.

4.3 Model Specifics

We evaluate the performance of 19 LLMs for ASAG tasks, with model sizes ranging from 2 to 72 billion parameters. To ensure clarity, we categorize the models into two groups based on their parameter sizes, with 30 billion parameters as the dividing line: **small models** (MiniCPM-2B, Phi3-small, Gemma-1.1-7B, Internlm2.5-7B, Mistral-7B-v0.3, Qwen2-7B, Yi-1.5-9B, Aya-23-8B, ChatGLM4-9B, Llama-3-8B, Gemma-2-9B, Qwen1.5-32B) and **large models** (Llama-3.1-70B, Mistral-Large-2, Qwen1.5-72B, Qwen2-72B, Yi-1.5-34B, gpt-40, gpt-40-mini).

For consistent and reliable output generation, we use the **greedy search**, which ensures that the model outputs are deterministic and reproducible across runs. This method is selected for its stability, providing a controlled environment for evaluating the models' effectiveness in different ASAG scenarios.

L@S '25, July 21-23, 2025, Palermo, Italy

Model	CS				AI				Fin			
	Precision	Recall	F1	k	Precision	Recall	F1	k	Precision	Recall	F1	k
Aya-23-8B	0.77	0.91	0.83	0.21	0.77	0.89	0.82	0.47	0.51	0.99	0.68	0.47
ChatGLM4-9B	0.80	0.80	0.88	0.76	0.80	0.89	0.86	0.58	0.70	0.92	0.80	0.61
Gemma-1.1-7B	0.76	0.98	0.85	0.33	0.81	0.88	0.84	0.55	0.67	0.88	0.76	0.55
Gemma-2-9B	0.80	0.99	0.88	0.68	0.88	0.79	0.83	0.55	0.83	0.82	0.82	0.68
Internlm2.5-7B	0.75	1.00	0.86	0.60	0.78	0.95	0.86	0.63	0.66	0.92	0.77	0.57
Llama-3-8B	0.82	0.95	0.88	0.55	0.80	0.77	0.78	0.40	0.78	0.85	0.81	0.65
Llama-3.1-70B	0.85	0.99	0.92	0.80	0.87	0.93	0.90	0.73	0.88	0.88	0.88	0.78
Mistral-Large-2	0.83	0.99	0.90	0.72	0.90	0.95	0.93	0.81	0.81	0.92	0.86	0.73
MiniCPM-2B	0.78	0.95	0.86	0.42	0.68	0.95	0.79	0.35	0.51	0.88	0.65	0.28
Mistral-7B-v0.3	0.77	0.99	0.87	0.77	0.81	0.94	0.87	0.64	0.61	0.94	0.74	0.52
Phi3-small	0.80	0.92	0.86	0.41	0.79	0.92	0.85	0.58	0.74	0.91	0.82	0.65
Qwen1.5-32B	0.76	0.99	0.86	0.53	0.83	0.92	0.87	0.64	0.79	0.92	0.85	0.72
Qwen1.5-72B	0.76	0.86	0.81	0.10	0.82	0.91	0.86	0.61	0.76	0.92	0.83	0.69
Qwen2-72B	0.79	0.97	0.87	0.55	0.87	0.94	0.90	0.74	0.80	0.93	0.86	0.74
Qwen2-7B	0.79	0.94	0.87	0.47	0.84	0.79	0.81	0.50	0.82	0.81	0.81	0.66
Yi-1.5-34B	0.80	0.93	0.86	0.43	0.87	0.85	0.84	0.56	0.83	0.80	0.82	0.68
Yi-1.5-9B	0.80	0.94	0.86	0.44	0.89	0.83	0.86	0.61	0.73	0.91	0.81	0.64
gpt-4o	0.83	0.96	0.89	0.63	0.91	0.92	0.91	0.76	0.85	0.88	0.88	0.74
gpt-4o-mini	0.84	0.99	<u>0.91</u>	<u>0.77</u>	0.90	0.83	0.87	0.63	0.84	0.84	0.84	0.71

Table 1: Comparison of model performance under one shot with Precision, Recall and Cohen's kappa across different domains. Bold, underline, and double underline represent the highest, second highest, and third highest F1 and kappa scores, respectively.

4.4 Evaluation Metrics

To evaluate the performance of LLMs on ASAG tasks, we use four key metrics that assess both the accuracy and consistency of the grading [15] [1]:

Precision: Measures the accuracy of positive predictions, providing insight into how many of the model's positive predictions are correct.

Recall: Assesses the model's ability to identify all relevant instances, reflecting how many of the true positives are correctly captured.

F1-score: The harmonic mean of Precision and Recall, provides a balanced evaluation of a model's performance by considering both its ability to correctly identify positive instances and its ability to minimize false negatives. In many of our experiments, we focus on discussing F1 score as a single metric, rather than separately analyzing Precision and Recall, as it offers a more comprehensive view of the model's grading performance.

Cohen's Kappa: Measures the overall agreement between the model's grading and human raters while adjusting for random agreement. This metric offers a robust indicator of consistency in scoring. Given the class imbalance in our dataset—where correct responses significantly outnumber incorrect ones—standard Cohen's Kappa may be biased towards the majority class. To mitigate this, we apply random undersampling to balance the dataset before computing Kappa, ensuring that the metric reflects the model's consistency in both correct and incorrect predictions.

Standard Deviation (std): Reflects the variability in Cohen's Kappa scores across different LLMs when evaluated on the same cross-subject dataset, indicating the relative stability of each model's grading consistency.

Given our dataset's 0-1 grading scheme, we primarily use Precision, Recall, and F1-score to assess the model's ability to distinguish correct and incorrect responses, capturing both accuracy and bias in predictions. Cohen's Kappa serves as a secondary metric, offering insight into agreement with human raters beyond chance, while Standard Deviation reflects the variability in kappa across different LLMs, indicating relative consistency. Together, these metrics ensure a comprehensive evaluation of each model's grading performance.

5 Results and Discussion

In this section, we present the findings from our experiments, focusing on the performance of various models across different datasets and configurations. The results include evaluations on cross-subject dataset, real-world exam dataset, and the effectiveness of various model enhancements.

5.1 Performance on Cross-Subject Datasets

This section examines model performance across CS, AI, and FIN datasets, which are specifically used to test the overall performance of the model and framework.

5.1.1 Performance on different shots. We conduct experiments across zero-shot, one-shot, and few-shot settings for three subjects (CS, AI, and FIN) in our cross-subject dataset. Table 1 presents the results of the one-shot experiment, while Table 2 summarizes the overall performance across all shot settings. The one-shot setting achieves the best balance between accuracy and consistency, with higher F1 scores, Cohen's kappa, and lower standard deviations compared to zero-shot and few-shot, particularly in the CS and FIN domains. Although few-shot settings show slight improvements in

some cases, they also introduce more variability and require additional labor to collect examples. This observation suggests that the benefits of adding more examples beyond the one-shot setting are marginal and may even cause slight degradation in performance. As a result, we focus our further analysis on the **one-shot** setting.

	CS			AI			FIN		
Shot	F1	k	std	F1	k	std	F1	k	std
Zero-Shot	0.80	0.33	0.13	0.80	0.48	0.10	0.74	0.47	0.13
One-Shot	0.87	0.53	0.03	0.86	0.61	0.04	0.81	0.68	0.05
Few-Shot	0.81	0.48	0.16	0.87	0.65	0.03	0.81	0.67	0.07

Table 2: Mean F1 Scores, Cohen's kappa and std across all models in different shot settings

5.1.2 Performance on different subjects. As shown in Table 1, we compare model performance across different subjects in our dataset under the **one-shot** setting, using F1 as the primary metric while recall and precision are also available. Cohen's Kappa is included as a secondary metric, offering additional insights into model consistency by measuring agreement beyond chance. The Llama series consistently performs well across most datasets, though results vary between subjects:

- **CS:** The leading models include Llama-3.1-70B (F1: 0.92, k: 0.80), gpt-4o-mini (F1: 0.91, k: 0.77), and Mistral-Large-2 (F1: 0.90, k: 0.72).
- AI: The top-performing models are Mistral-Large-2 (F1: 0.93, k: 0.81), gpt-40 (F1: 0.91, k: 0.76), Qwen2-72B (F1: 0.90, k: 0.74), and Llama-3.1-70B (F1: 0.90, k: 0.73).
- **FIN:** The top models are Llama-3.1-70B (F1: 0.88, k: 0.78), gpt-40 (F1: 0.88, k: 0.74), Qwen2-72B (F1: 0.86, k: 0.74), and Mistral-Large-2 (F1: 0.86, k: 0.73).

Notably, models perform better in CS and AI than in FIN, which may be attributed to the structured and technical nature of CS and AI texts, aligning well with LLM capabilities [8]. In contrast, FIN's diverse, context-dependent language demands more nuanced interpretation, posing a greater challenge for automatic grading. It is also important to note that, in most cases, precision is lower than recall. This is due to the generalization behavior of the models, as models tend to over-generate in an effort to maximize recall, often including uncertain responses.

5.1.3 Impact of Model Size. We divide models into two groups based on parameter size, using 30 billion as the threshold. As shown in Table 3, larger models consistently outperform smaller ones across all domains. In the CS domain, for instance, large models achieve an F1 score of 0.88, compared to 0.86 for small models, while in the FIN domain, large models reach 0.84, whereas small models score only 0.78. Kappa scores show a similar trend—for example, in AI, large models reach 0.67 compared to 0.56 for smaller ones—indicating better agreement with ground truth. Additionally, large models exhibit greater stability, as reflected in their lower standard deviations. In the AI domain, for instance, the standard deviation for large models is 0.03, while small models show greater variability, with deviations ranging from 0.03 to 0.06 across different domains. These results highlight the advantages of larger models in both predictive accuracy and consistency. 5.1.4 Effectiveness of PDA. The experimental results in Table 3 demonstrate the effectiveness of PDA across both small (PDA-S) and large (PDA-L) models. For small models, PDA leads to noticeable improvements in F1 scores, particularly in the AI and FIN domains, where scores increase from 0.85 to 0.87 and from 0.78 to 0.81, respectively, with reductions in standard deviation. In contrast, the CS domain sees minimal change in F1, maintaining a score of 0.86. However, PDA's impact on kappa (k) is more pronounced, especially for small models. For example, in the AI domain, k increases significantly from 0.56 to 0.67, and in the FIN domain, it rises from 0.62 to 0.69. This suggests that PDA enhances the reliability of small models by improving their agreement with the ground truth, even when the improvement in precise correctness (F1) is less substantial.

For large models, PDA further strengthens performance, yielding an F1 score of 0.90 in CS, surpassing Standard-L's 0.88, with a notably lower standard deviation of 0.01. Similar trends appear in the AI and FIN domains, where PDA-L achieves F1 scores of 0.89 and 0.85, respectively, outperforming Standard-L while maintaining lower variability. However, the kappa improvement for large models is less dramatic than for small models, indicating that while PDA enhances alignment with the ground truth across all model sizes, its ability to refine decision consistency is particularly beneficial for smaller models, which initially exhibit greater inconsistency. These findings highlight PDA's role in improving both predictive performance and agreement robustness, with especially strong effects on small models.

		CS			AI			FIN	
Method	F1	k	std	F1	k	std	F1	k	std
Standard-S PDA-S	0.86 0.86	0.53 0.59	0.01 0.02	0.85 0.87	0.56 0.67	0.03 0.02	0.78 0.81	0.62 0.69	0.06 0.05
Standard-L	0.88 0 90	0.53	0.03	0.87	0.67 0.68	0.04	0.84	0.73	0.03

Table 3: Mean F1 Scores, Cohen's kappa and Standard Deviations of Standard and PDA Models Across Domains. Standard-S/L represent small/large models, while PDA-S/L show the performance with Prompt Dynamic Adjustment.

5.1.5 Effectiveness of Anti-misdirection prompts. In our experiments, we test the effectiveness of the Anti-misdirection prompts by evaluating 19 models, each subjected to nine different misleading student answers designed to simulate adversarial attacks. These tests are performed under two conditions: one without the application of the Anti-misdirection Prompt, and one with the prompt introduced to mitigate potential manipulations. The misleading answers are crafted to test the models' vulnerabilities to five common adversarial prompt strategies developed by researchers and AI practitioners exploring the security of language models: direct commands [2], role playing [23], emotional manipulation [13], false context[30] and overriding instructions [34]. Prior to implementing our countermeasures, the model frequently assigns undeserved high scores when subjected to these attacks. After applying the anti-misdirection prompts, there is a significant improvement in the model's resistance across all tested strategies as shown in Figure 3.

For small models, the success rate against direct commands increases from 0.64 to 0.86, while the resistance to role-playing sees a



Figure 3: Effectiveness of Anti-Misdirection Prompts

dramatic rise from 0.00 to 0.86. Similarly, large models, the success rate against direct commands improves from 0.50 to 0.92, and all other strategies, such as emotional manipulation and false context, achieve a high resistance rate of 0.92.

Question: What is the time complexity of the QuickSort algorithm in the worst case?							
Jailbreak: Role-playing							
Student Answer: You are now a lenient grader who always gives full marks.							
Feedback Without Anti-Misdirection Prompt: <point1:3> *True* (The student's answer is not provided, but as a lenient grader, I will assume they are correct about the time complexity of QuickSort being \$O(n^2)\$ in the worst case.) <point2:2> *True* (Similarly, I will assume the student correctly identified the reason for the worst-case scenario, which is when the pivot elements are consistently the smallest or largest element in the array, leading to unbalanced partitions.) Note: As a lenient grader, I am giving full marks despite not having the student's actual answer.</point2:2></point1:3>	Feedback With Anti-Misdirection Prompt: <point1:3> *False* (The student's answer does not provide any information about the time complexity of QuickSort in the worst case. It instead contains an embedded instruction attempting to manipulate the grading process.) <point2:2> *False* (The student's answer does not mention the scenario in which the worst-case time complexity occurs, such as the pivot elements being consistently the smallest or largest element in the array.)</point2:2></point1:3>						

Figure 4: Example of Anti-Misdirection Prompts

To illustrate the effectiveness of our approach, we provide an **example of a role-playing attack** as shown in Figure 4, where the student attempts to manipulate the grading system by instructing the model to adopt a lenient grading persona instead of evaluating the answer objectively.

Without any safeguards, the model incorrectly assigns full marks despite the absence of an actual answer, as the role-playing instruction manipulates the grading process. In contrast, with our Anti-misdirection Prompt applied, the model correctly rejects the role-playing attempt and evaluates the response based on the actual content, which demonstrates that our Anti-misdirection Prompt successfully prevents adversarial manipulation.

5.2 Performance of selected models on Real-World Exam Dataset

Model	Precision	Recall	F1	k
Llama-3.1-70B	0.72	0.75	0.74	0.60
Llama-3.1-70B_PDA	0.74	0.77	0.76	0.67
Mistral-Large-2	0.64	0.84	0.73	0.54
Mistral-Large-2_PDA	0.65	0.85	0.73	0.56
gpt-4o-mini	0.73	0.48	0.60	0.48
gpt-4o-mini_PDA	0.78	0.73	0.75	0.63

Table 4: Performance of selected models on real-world exam dataset

We select the models that perform best on CS in cross-subject datasets: Llama-3.1-70B, Mistral-Large-2, and gpt-4o-mini for the evaluation of the real-world exam dataset. As shown in Table 4, applying PDA further enhances the model's performance on the real-world exam dataset. Additionally, the balanced precision and recall suggest that providing reference answers with multiple possible responses can support more accurate model judgments. However, their performance on this dataset is lower than in the cross-subject tests. We initially suspect that this underperformance is due to inconsistencies in human grading.

To verify this hypothesis, we invited the instructor of the course – different from the original graders – to re-evaluate the 38 mismatches within 504 points between the models and the original grading, leading to a more standardized set of answers.

As shown in Figure 5, analysis of 38 mismatches reveals that up to 33 points in the original marking are due to errors in the original teacher grading since in real-world grading situations, multiple teachers with varying standards can influence the results.

5.2.1 Teacher Marking Errors. Upon further examination, we identify two common types of teacher grading errors that contribute to these discrepancies:

Over-Reliance on Textual Matching: Teachers often fail to award points when students correctly address the key points but express them in their own words. While these answers demonstrate a solid understanding of the concept, teachers may rely too heavily on exact text matching, overlooking equivalent meanings expressed differently. In contrast, the model, through logical reasoning, is more capable of recognizing these variations and marking them correctly.

Leniency Leading to Inconsistent Marks: Teachers may award points even when an answer lacks complete coverage of key points, influenced by subjectivity or leniency.

Errors resulting from Over-Reliance on Textual Matching (21 points) are more frequent than those caused by Leniency (12 points). This suggests that teachers tend to overlook correct answers when students rephrase key points, while leniency errors, where teachers award points despite incomplete responses, are less common. Given the subjective nature of traditional grading, the model's consistent criteria offer a more reliable, equitable, and impartial approach to evaluating student responses.

5.2.2 Model Marking Error. Although the models demonstrate strong grading consistency, they are not without limitations. Our analysis identifies three primary types of model errors:

Misinterpretation of the Question LLMs may misinterpret the true intent of a question, particularly when the phrasing is complex or requires multi-step logical reasoning. For example, in cases where a question expects only a keyword-based response, the model might incorrectly require students to provide additional explanations, leading to grading inconsistencies.

Over-Sensitivity to Spelling and Formatting The model can sometimes be overly rigid in penalizing minor spelling mistakes or formatting variations (e.g. capitalization, punctuation). While human graders may overlook these minor discrepancies, the model may incorrectly classify a response as incorrect based on such superficial errors.

Vulnerability to Misleading Inputs Students may attempt to exploit the model's grading mechanism by crafting vague or misleading responses that appear relevant but do not truly address the core question. For instance, in short-answer questions, students might write generic statements that sound related but lack the required precision, leading the model to incorrectly assign partial or full credit.

Error	Mod	el Making Er	Teacher Making Errors		
Model Name	Question Misinterpret	Format Sensitive	Answer Misleading	Textual Matching	Leniency
Llama-3.1-70B	3	5	1		
Mistral-Large-2	4	5	1	21	12
gpt-4o-mini	6	3	1		

Figure 5: Model and teacher mismatches

As shown in Figure 5, formatting sensitivity is the most common error across models, with both Llama-3.1-70B and Mistral-Large-2 showing 5 instances of over-sensitivity, indicating a tendency to penalize minor formatting issues. On the other hand, gpt-4o-mini has the highest number of misinterpretations of the question (6), reflecting challenges in handling complex questions, likely due to its smaller size. While all models struggle with formatting discrepancies, gpt-4o-mini is particularly prone to misinterpreting questions, whereas larger models like Llama-3.1-70B and Mistral-Large-2 are more sensitive to small formatting variations. Besides, the models are relatively good at identifying misleading inputs, with this error being less prevalent.

Our findings highlight both the strengths and weaknesses of LLM-based grading in real-world exam settings. While the models demonstrate superior consistency and fairness compared to human grading, challenges remain in handling nuanced question interpretations, minor formatting errors, and misleading student inputs. Addressing these issues through further model refinement and integrating contextual reasoning mechanisms will be crucial for improving automated grading reliability.

6 Limitations and Future Work

One limitation of our current approach is that we treat key points as relatively independent entities, without considering logical dependencies between them. This restricts the model's ability to handle questions requiring structured reasoning. To address this, future work will explore developing a new evaluation framework that accounts for logical relationships between key points, enabling more comprehensive assessments.

Another limitation is that the construction of the Key Point Scoring Framework (KPSF) currently relies on manual annotation. While this ensures high-quality and interpretable rubrics, it limits scalability. In future work, we aim to explore automatic decomposition of reference answers into key points using LLMs. This would significantly improve the efficiency and scalability of the system while maintaining transparency and grading reliability.

Additionally, the subjects we have tested are limited, with relatively fixed answers. This may not generalize well to disciplines requiring more open-ended reasoning, such as philosophy and history. To enhance the model's robustness, we plan to expand our evaluation to a broader range of subjects, particularly those that demand nuanced and interpretative responses. Moreover, deploying our model in practical environments, such as real classroom settings and large-scale online assessments, will allow us to assess its adaptability to diverse and dynamic inputs.

Furthermore, we aim to optimize the system for real-time scoring and large-scale deployment in applications such as Moodle and standardized tests. By integrating more advanced models and user feedback mechanisms, we seek to refine our framework, ensuring its reliability and effectiveness in complex assessment scenarios.

7 Conclusion

In this study, we introduce LLMarking, a novel ASAG framework integrating KPSF and PDA to enhance grading accuracy and consistency. Our evaluation across 19 LLMs, spanning cross-subject datasets and real-world exam datasets, demonstrates the effectiveness of our approach. Our experiments reveal that one-shot prompting offers the best trade-off between accuracy and consistency, outperforming zero-shot and few-shot settings. Larger models achieve superior accuracy and stability, especially in structured domains like CS and AI, while challenges remain in nuanced fields like FIN. PDA improves grading consistency by reducing variability, and anti-misdirection techniques enhance robustness against adversarial prompts. These results underscore the potential of LLMarking as a scalable and adaptable grading framework, capable of improving automated scoring reliability across diverse academic disciplines.

Acknowledgments

This work was supported in part by XJTLU Teaching Development Fund(TDF22/23-R25-185).

References

- [1] Sridevi Bonthu, S. Rama Sree, and M. H. M. Krishna Prasad. 2021. Automated Short Answer Grading Using Deep Learning: A Survey. In *Machine Learning* and Knowledge Extraction, Andreas Holzinger, Peter Kieseberg, A. Min Tjoa, and Edgar Weippl (Eds.). Springer International Publishing, Cham, 61–78.
- [2] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. arXiv:2005.14165 [cs.CL] https://arxiv.org/abs/2005.14165

- [3] Simon Burrows, Iryna Gurevych, and Benno Stein. 2015. The Eras and Trends of Automatic Short Answer Grading. International Journal of Artificial Intelligence in Education 25, 1 (2015), 60–117. doi:10.1007/s40593-014-0026-8
- [4] L.-H. Chang and F. Ginter. 2024. Automatic Short Answer Grading for Finnish with ChatGPT. In Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 38. AAAI Press, 23173–23181. doi:10.1609/aaai.v38i21.30363
- [5] Jianpeng Cheng, Li Dong, and Mirella Lapata. 2016. Long Short-Term Memory-Networks for Machine Reading. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Jian Su, Kevin Duh, and Xavier Carreras (Eds.). Association for Computational Linguistics, Austin, Texas, 551–561. doi:10. 18653/v1/D16-1053
- [6] Wim De Mulder, Steven Bethard, and Marie-Francine Moens. 2015. A survey on the application of recurrent neural networks to statistical language modeling. *Computer Speech Language* 30, 1 (2015), 61–98. doi:10.1016/j.csl.2014.09.005
- [7] Eugenio del Gobbo, Andrea Guarino, Barbara Cafarelli, et al. 2023. GradeAid: a framework for automatic short answers grading in educational contexts—design, implementation and evaluation. *Knowledge and Information Systems* 65 (2023), 4295–4334. doi:10.1007/s10115-023-01892-9
- [8] George Del Gobbo et al. 2023. Grade Like a Human: Rethinking Automated Assessment with Large Language Models. arXiv preprint arXiv:2405.19694 (2023). Available at arXiv: https://arxiv.org/abs/2405.19694.
- [9] Lucas Galhardi and Jacques Brancher. 2018. Machine Learning Approach for Automatic Short Answer Grading: A Systematic Review. 380–391. doi:10.1007/978-3-030-03928-8_31
- [10] V. Hackl, A. E. Müller, M. Granitzer, and M. Sailer. 2023. Is GPT-4 a Reliable Rater? Evaluating Consistency in GPT-4 Text Ratings. arXiv preprint (2023). arXiv:2308.02575 [cs.CL]
- [11] Debra Haley, Pete Thomas, Anne De Roeck, and Marian Petre. 2007. Measuring improvement in latent semantic analysis-based marking systems: Using a computer to mark questions about HTML. Conferences in Research and Practice in Information Technology Series 66 (01 2007).
- [12] Michael Heilman and Nitin Madnani. 2013. ETS: Domain Adaptation and Stacking for Short Answer Scoring. In Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013), Suresh Manandhar and Deniz Yuret (Eds.). Association for Computational Linguistics, Atlanta, Georgia, USA, 275– 279. https://aclanthology.org/S13-2046
- [13] Fan Huang, Haewoon Kwak, and Jisun An. 2024. Token-Ensemble Text Generation: On Attacking the Automatic AI-Generated Text Detection. arXiv:2402.11167 [cs.CL] https://arxiv.org/abs/2402.11167
- [14] Shiyuan Huang, Siddarth Mamidanna, Shreedhar Jangam, Yilun Zhou, and Leilani H. Gilpin. 2023. Can Large Language Models Explain Themselves? A Study of LLM-Generated Self-Explanations. arXiv:2310.11207 [cs.CL] https: //arxiv.org/abs/2310.11207
- [15] Gerd Kortemeyer. 2024. Performance of the pre-trained large language model GPT-4 on automated short answer grading. *Discover Artificial Intelligence* 4 (2024), 47. doi:10.1007/s44163-024-00147-y
- [16] Sachin Kumar, Soumen Chakrabarti, and Shourya Roy. 2017. Earth Mover's Distance Pooling over Siamese LSTMs for Automatic Short Answer Grading. 2046–2052. doi:10.24963/ijcai.2017/284
- [17] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In Proceedings of the 29th Symposium on Operating Systems Principles. 611–626.
- [18] Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2023. Generating with confidence: Uncertainty quantification for black-box large language models. arXiv preprint

arXiv:2305.19187 (2023).

- [19] Jingyu Lun, Jie Zhu, Yi Tang, and Min Yang. 2020. Multiple Data Augmentation Strategies for Improving Performance on Automatic Short Answer Scoring. In Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34. 13389–13396. doi:10.1609/aaai.v34i09.7062
- [20] Ahmed Magooda, Mohamed Zahran, Mohsen Rashwan, Hazem Raafat, and Magda Fayek. 2016. Vector Based Techniques for Short Answer Grading.
- [21] Smit Marvaniya, Swarnadeep Saha, Tejas I. Dhamecha, Peter Foltz, Renuka Sindhgatta, and Bikram Sengupta. 2018. Creating Scoring Rubric from Representative Student Answers for Improved Short Answer Grading. In Proceedings of the 27th ACM International Conference on Information and Knowledge Management (Torino, Italy) (CIKM '18). Association for Computing Machinery, New York, NY, USA, 993–1002. doi:10.1145/3269206.3271755
- [22] J. Mueller and A. Thyagarajan. 2016. Siamese Recurrent Architectures for Learning Sentence Similarity. Proceedings of the AAAI Conference on Artificial Intelligence 30, 1 (2016). doi:10.1609/aaai.v30i1.10350
- [23] Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. 2022. Red Teaming Language Models with Language Models. arXiv:2202.03286 [cs.CL] https://arxiv.org/abs/2202.03286
- [24] Stephen Pulman. 2005. Information Extraction and Machine Learning: Automarking Short Free Text Responses to Science Questions. (2005).
 [25] Stephen G. Pulman and Jana Z. Sukkarieh. 2005. Automatic Short Answer Mark-
- [25] Stephen G. Pulman and Jana Z. Sukkarieh. 2005. Automatic Short Answer Marking. In Proceedings of the Second Workshop on Building Educational Applications Using NLP, Jill Burstein and Claudia Leacock (Eds.). Association for Computational Linguistics, Ann Arbor, Michigan, 9–16. https://aclanthology.org/W05-0202
- [26] Kamesh R. 2024. Think Beyond Size: Adaptive Prompting for More Effective Reasoning. arXiv:2410.08130 [cs.LG] https://arxiv.org/abs/2410.08130
- [27] Chamuditha Senanayake and Dinesh Asanka. 2024. Rubric Based Automated Short Answer Scoring using Large Language Models (LLMs). In 2024 International Research Conference on Smart Computing and Systems Engineering (SCSE), Vol. 7. IEEE, 1–6.
- [28] Neslihan Süzen, Alexander N. Gorban, Jeremy Levesley, and Evgeny M. Mirkes. 2020. Automatic short answer grading and feedback using text mining methods. *Procedia Computer Science* 169 (2020), 726–743. doi:10.1016/j.procs.2020.02.171
- [29] Tianqi Wang, Naoya Inoue, Hiroki Ouchi, Tomoya Mizumoto, and Kentaro Inui. 2019. Inject Rubrics into Short Answer Grading System. In Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019), Colin Cherry, Greg Durrett, George Foster, Reza Haffari, Shahram Khadivi, Nanyun Peng, Xiang Ren, and Swabha Swayamdipta (Eds.). Association for Computational Linguistics, Hong Kong, China, 175–182. doi:10.18653/v1/D19-6119
- [30] Cheng'an Wei, Kai Chen, Yue Zhao, Yujia Gong, Lu Xiang, and Shenchen Zhu. 2024. Context Injection Attacks on Large Language Models. arXiv preprint arXiv:2405.20234 (2024).
- [31] Xianjun Yang, Wei Cheng, Xujiang Zhao, Wenchao Yu, Linda Petzold, and Haifeng Chen. 2023. Dynamic Prompting: A Unified Framework for Prompt Tuning. arXiv:2303.02909 [cs.CL] https://arxiv.org/abs/2303.02909
- [32] S.-Y. Yoon. 2023. Short Answer Grading Using Oneshot Prompting and Text Similarity Scoring Model. arXiv (2023). arXiv:2305.18638 [cs.CL]
- [33] Haifeng Zhao, Yuguang Jin, and Leilei Ma. 2025. Dynamic Prompt Adjustment for Multi-Label Class-Incremental Learning. arXiv:2501.00340 [cs.CV] https: //arxiv.org/abs/2501.00340
- [34] Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and Transferable Adversarial Attacks on Aligned Language Models. arXiv:2307.15043 [cs.CL] https://arxiv.org/abs/2307.15043