# The Evolution of Short Answer Grading Systems: From Manual Methods to AI-Driven Solutions with GPT-4

### Augustine O. Ugbari
Department of Computer Sciences
University of Port Harcourt
Choba, Nigeria

### Chidiebere Ugwu
Department of Computer Sciences
University of Port Harcourt
Choba, Nigeria

### Laeticia N. Onyejegbu
Department of Computer Sciences
University of Port Harcourt
Choba, Nigeria

## ABSTRACT
The evaluation of short answer responses has long been a critical component of educational assessments, providing insights into student comprehension and analytical skills. This paper traces the evolution of short answer grading systems from their inception in manual grading practices to the advent of AI-driven solutions, focusing particularly on the advancements brought by models like GPT-4. Through a comprehensive review of historical developments, technological advancements, and pedagogical impacts, this research provides a detailed understanding of the progression and future prospects of short answer grading systems.

## Keywords
Short Answer Grading System, GPT-4, Machine Learning.

## 1. INTRODUCTION
Over the past few years, the rapid expansion of educational technology has played a significant role in shaping the modern learning landscape. As online learning platforms and massive open online courses (MOOCs) have become increasingly popular, the demand for effective and reliable automated assessment systems has grown substantially. Short answer questions (SAQs) has proven to an essential tool in education, testing students' knowledge, understanding, and ability to articulate responses concisely. The grading of these responses has traditionally been a labor-intensive process, demanding significant time and effort from educators. (Blum et al., 2020)

The process of grading short answer questions manually involves an evaluator reading through each student response and assessing its quality based on a set of predefined criteria. The criteria for grading may vary depending on the specific question being asked and the educational context, but they generally include factors such as relevance, accuracy, completeness, and coherence (Rodgers & Beeson, 2009). To begin the process, the evaluator first reads through each response to get a sense of its overall quality and content. They then evaluate each response against the criteria for grading, assigning a score or grade based on how well the response meets each criterion. The evaluator may also provide comments or feedback to the student to help them understand why they received a particular grade and how they can improve their future responses.

Grading short answers manually can be both time-consuming and subjective, as it heavily relies on the evaluator's discretion and understanding of the grading criteria. This subjectivity often results in inconsistencies, as different evaluators may assign disparate grades to identical responses. To mitigate these issues, automated short answer grading systems (ASAG) have been developed, offering a more objective and streamlined grading process. Leveraging machine learning algorithms and natural language processing techniques, these systems assess student responses against predefined criteria, thereby ensuring a fairer and more efficient grading process.

## 2. METHODOLOGY
This study employs a multi-faceted research approach, including literature reviews, case studies, technological analysis, comparative studies, and expert interviews. Each method provides unique insights into different aspects of the evolution of short answer grading systems.

## 3. LITERATURE REVIEW
### 3.1 Origins and Early Practices
The grading of short answer responses has its roots in the earliest forms of formal education. Initially, grading was entirely manual, with teachers evaluating each student's response based on predefined criteria. This process, while thorough, was time-consuming and prone to subjective biases.

Manual grading has faced several significant challenges. Firstly, subjectivity is a major issue, as different teachers might grade the same response differently, leading to inconsistencies. Additionally, even the same teacher can grade inconsistently at different times due to factors like fatigue or mood, which compromises the fairness and reliability of assessments.

Manual grading is also very time-consuming, requiring significant effort to evaluate a large number of responses. This makes scalability a problem, especially as class sizes grow, making manual grading impractical for large-scale assessments. The time needed for manual grading often results in delayed feedback to students, which is detrimental since timely feedback is crucial for effective learning and improvement.

Furthermore, the effort involved in manual grading can limit the depth and specificity of feedback provided. Teachers might resort to generic comments or simply assign a grade without detailed explanations, reducing the instructional value of their feedback. Lastly, grading short answer questions is labor-intensive, consuming substantial time and effort from teachers, which detracts from the time they could spend on instructional activities and professional development.

### 3.2 This paragraph is a repeat of 3.1
Please use a 9-point Times Roman font, or other Roman font with serifs, as close as possible in appearance to Times Roman in which these guidelines have been set. The goal is to have a 9-point text, as you see here. Please use sans-serif or non-proportional fonts only for special purposes, such as distinguishing source code text. If Times Roman is not available, try the font named Computer Modern Roman. On a Macintosh, use the font named Times. Right margins should be justified, not ragged.

## 3.3 Early Innovations and Attempts at Automation

One of the earliest attempts to automate grading involved Optical Mark Recognition (OMR) technology, which emerged in the mid-20th century. While OMR was primarily used for multiple-choice questions, it laid the groundwork for exploring automation in other types of assessments. OMR systems could quickly and accurately process large volumes of answer sheets, demonstrating the potential for technology to streamline grading processes. One of the pioneers in this field was Michael Sokolski, who co-founded the company Scantron Corporation in 1972, which would become synonymous with OMR technology.

Key Research Milestones on OMR include Early Optical Scanners (1960s) (IBM Archives, 2023). The Statistical Machine could read and interpret data from cards punched with specific patterns, which were manually created (Riley, 1955). Early Optical Scanners (Roberts, 1963) was used for the grading of standardized tests, where students filled out answer sheets by marking bubbles corresponding to their answers. Michael Sokolski and Scantron (Sokolski, 1972), the Scantron's machines used optical sensors to detect the presence of marks on specially designed forms.

Despite its advantages, early OMR technology faced several limitations such as their Form Design, Mark Quality and Initial Cost.

In the latter half of the 20th century, researchers began experimenting with simple computational models for grading short answer responses. These early models were rudimentary and typically focused on basic pattern matching and keyword recognition.

Keyword Matching introduced the concept of text similarity. These systems compared student responses against a list of predefined keywords. If the response contained the required keywords, it was considered correct. While this approach was straightforward, it lacked the ability to understand context or nuance in student responses. Various forms of Keyword Matching exist, such as Exact Keyword Matching (Foltz et al., 1999), Partial Keyword Matching (Jurafsky & Martin, 2020), Synonym-Based Keyword Matching (Manning et al., 2008), Stemming and Lemmatization (Wilcox, 2020).

Pattern Recognition: Some early systems attempted to recognize patterns or phrases within responses. These systems were slightly more sophisticated than keyword matching but still limited in their ability to handle variations in student language and expression. (Duda et al., 2001)

There are various implementations of Pattern Recognition such as:

1. Template matching is one approach, involving the comparison of input patterns with stored templates to find the closest match (Bishop, 2006).

2. Statistical pattern recognition uses statistical techniques to model and classify patterns based on their features and distributions (Goodfellow et al., 2016).

3. Syntactic, or structural, pattern recognition focuses on the relationships between components of patterns, using grammatical rules to describe and recognize complex structures (Fu, 1982).

4. Neural network-based pattern recognition leverages the capabilities of neural networks to learn and identify patterns

through training on large datasets, improving accuracy and adaptability over time (Goodfellow et al., 2016).

## 3.4 Similarity Measures

The concept of similarity serves as a broad umbrella term encompassing a diverse array of scores and measures designed to evaluate distinctions and relationships within various types of data. After gathering data, one of the initial inquiries that often captivates researchers is the extent of similarity between two data samples, be it a text, an individual, or an event. This question holds considerable significance across diverse fields, particularly in the realms of humanities and critical analysis, where the comparison of similarities and disparities between two entities remains a fundamental pursuit. It's important to note that beyond computational methods, non-computational assessments of similarities and differences underpin a substantial portion of critical endeavours and scholarly activities. For example, the genre of a text may be determined by considering its similarity to similar texts which have so far been identified as being part of that category. On the other hand, new sources of criticism could be opening when it is known that a certain text differs widely from another in an established genre. In academic categorization and critical analysis, the uniqueness of a study object or its similarity in relation to others or groups may play an important role.

Statistical similarity indicators form the foundation for various clustering and classification methods, enabling researchers to quantify the similarity or dissimilarity between objects under study. In text analysis, one can evaluate the similarity of two texts by representing each text as a sequence of word counts and computing the distance based on these word counts as features.

Four various approaches may be used to classify similar measures. String Based Similarity, Corpus Based Similarity, Knowledge Based Similarity and Hybrid Similarity Measures.

String similarity measures analyse string sequences and the composition of characters to determine their similarity. These measures play a crucial role in various applications, including natural language processing, data deduplication, and fuzzy string matching. These measures are encapsulated within what is known as a "string metric," which quantifies the degree of similarity or dissimilarity (often referred to as distance) between two strings of text. These metrics are utilised in matching strings approximately which aids in tasks where assessing the similarity or variance between strings is paramount. They encompass both character-based (LCS, Damerau Levenshtein, Jaro, Jaro-Winkler, Needleman-Wunsch, N-gram) and term-based approaches (Block Distance, Cosine Similarity, Dice's Coefficient, Euclidean Distance, Jaccard similarity, Matching Coefficient, Overlap Coefficient) all of which offers a range of methods for assessing string comparison.

In Corpus-based similarity, these methods are used to quantify the similarity between pieces of text based on the information derived from a large corpus of text data. Instead of focusing solely on the presence or absence of specific words or characters, corpus-based similarity takes into account the frequency of occurrence of words, phrases, or other linguistic units in a given text corpus. This approach leverages statistical analysis and natural language processing techniques to compute similarity scores between texts. By analyzing patterns of word usage and co-occurrence in a corpus, corpus-based similarity methods can capture semantic relationships and contextual information that contribute to the overall similarity between texts.

Knowledge-Based Similarity represents one of the semantic similarity measurement methods that rely on assessing the level of similarity between words through insights drawn from semantic networks. It refers to the assessment of similarity between terms or concepts based on external knowledge sources, such as ontologies, semantic networks, or knowledge bases such as WordNet. Unlike distributional or corpus-based methods that rely solely on statistical patterns in text data, knowledge-based similarity considers the semantic relationships and hierarchical structures present in external knowledge sources. WordNet (Miller et al., 1990) stands out as the most widely used semantic network for quantifying Knowledge-Based similarity between words. WordNet serves as an extensive lexical database for the English language. Within WordNet, nouns, verbs, adjectives, and adverbs are organized into sets called synsets. Each synset represents a distinct cognitive concept, allowing users to explore the network of interconnected words and their meanings. These synsets are interconnected through a network of conceptual-semantic and lexical relationships, providing a rich resource for semantic similarity assessments.

Hybrid similarity measures is a method that combines more than one similarity measures to enhance the accuracy and robustness of semantic similarity assessments. By leveraging the strengths of different approaches, these hybrid techniques provide a more comprehensive understanding of relatedness between concepts. This can be a combination of semantic similarity measures, corpus-based measures, and knowledgebase. Each class of similarity measure can consist of multiple type, such that we can evaluate different semantic algorithms separately, before then combining their several similarity metrics into one..

# 4. TECHNOLOGICAL PROGRESSION: AUTOMATED GRADING SYSTEMS

## 4.1 Rule-Based Systems

The first significant step towards automation in grading short answers was the development of rule-based systems. These systems used predefined rules to evaluate responses. Rule-based systems are automated systems that use a set of predefined rules to evaluate and grade student responses. These rules are typically created by experts and are based on the expected answers to the questions. The system utilizes these rules to evaluate student responses, determining their accuracy and assigning grades accordingly [1]. Its characteristics include explicit rules, which are clearly defined and encoded into the system, ensuring transparency in the grading process. Moreover, the system operates deterministically, consistently producing the same output for the same input based on the predefined rules, further enhancing its transparency and reliability [2].

Operation involves several key processes in the grading system [3]. Firstly, experts establish rules dictating the criteria for a correct answer. These rules encompass factors like keyword matching, specific phrases, and logical structures. Secondly, student responses undergo analysis based on these rules. The system evaluates whether required keywords, phrases, and logical structures are present in the answer. Finally, grading occurs where the system assigns a grade to the response based on the analysis conducted. If the response aligns with the predefined rules, it receives full or partial credit.

In an illustrative scenario, envision a short answer query inquiring about the primary causes of the American Civil War. Within a rule-based system, specific guidelines are established:

The response is expected to contain essential keywords such as "slavery," "states' rights," and "economic differences." Additionally, responses that incorporate "slavery" alongside at least one additional keyword may qualify for partial credit.

In various contexts, rule-based systems find application: In standardized testing scenarios, these systems are utilized to evaluate short answer and essay questions, particularly when expected responses can be distinctly defined.

Similarly, in formative assessments within classrooms, teachers rely on rule-based systems to promptly assess students' comprehension of specific concepts, providing immediate feedback.

These systems offer several advantages:

1. Consistency: By uniformly applying criteria to all responses, rule-based systems ensure consistent grading outcomes.

2. Efficiency: They possess the capability to grade responses at a significantly faster pace than human graders, facilitating swift turnaround times.

3. Transparency: Utilizing explicit rules renders the grading process transparent and easily comprehensible.

However, they also present certain limitations:

1. Rigidity: Rule-based systems may encounter challenges in accommodating variations in responses, especially those that are correct but articulated differently from the predefined rules.
2. Scalability: The process of creating and maintaining a comprehensive set of rules for intricate subjects can be time-consuming and daunting.
3. Lack of Contextual Understanding: These systems lack the capability to grasp context and nuance, potentially resulting in the mis-grading of complex or ambiguous answers.

Early rule-based grading systems were simple and relied heavily on keyword matching and basic logical rules. These systems were primarily used in contexts where the expected answers were straightforward and easy to predict. Take for instance, an early system might be used to grade factual recall questions in history, such as "Who was the first president of the United States?" with the rule being the response must contain "George Washington."

Over time, rule-based systems evolved to incorporate more sophisticated rules and logic. This included:

(a) **Advanced Keyword Matching**: Incorporating partial and synonym-based keyword matching to handle variations in student responses.

(b) **Contextual Rules**: Developing rules that consider the context and structure of the response, not just the presence of specific keywords.

(c) **Hierarchical Rules**: Using hierarchical and nested rules to evaluate more complex responses that involve multiple concepts and ideas.

In a more advanced system, a question asking for an explanation of photosynthesis might involve rules that check for a sequence of keywords ("light," "chlorophyll," "carbon dioxide," "glucose") and their relationships (e.g., "light is absorbed by chlorophyll").

## 4.2 Machine Learning Models

The advent of machine learning introduced more sophisticated approaches to automated grading. These models could learn from examples and improve their accuracy over time. [18]

Machine learning models are algorithms that learn patterns from data to make predictions or decisions without being explicitly programmed. In the context of automated grading, these models are trained on a dataset of graded responses and learn to predict grades for new, unseen responses based on the patterns identified during training.

In terms of functionality, machine learning models underwent two critical processes:

1. Training Data: These models were trained on extensive datasets comprising graded responses, allowing them to learn from a diverse range of examples.

2. Pattern Recognition: Leveraging their capabilities, these models could discern intricate patterns and relationships within the data, enabling more nuanced analysis and grading.

In terms of advantages, machine learning models demonstrated greater flexibility and accuracy compared to rule-based systems. They exhibited the capability to accommodate a broader array of responses and refine their performance over time. Conversely, these models faced drawbacks such as their dependence on extensive training data and their computational intensity. Characteristics of these systems include their reliance on extensive datasets to discern grading patterns. They exhibit adaptability, capable of refining their performance with the accumulation of additional data. However, they operate in a non-deterministic manner, potentially yielding varying outcomes for similar inputs based on learned patterns [1].

Operation involves a series of steps [4]: Firstly, data collection entails gathering a substantial dataset comprising student responses paired with their respective grades. Next, feature extraction is performed to identify critical attributes from the responses. These features encompass various elements such as keywords, sentence structure, semantic meaning, among others. Subsequently, model training takes place, wherein the extracted features and grades are utilized to train the machine learning model. Common algorithms employed for this purpose include decision trees, support vector machines (SVM), and neural networks. Finally, prediction occurs once the model is adequately trained. This involves the model's capability to forecast grades for new student responses by analyzing their features and comparing them to the learned patterns.

For a question on explaining the process of evaporation, a machine learning model might be trained on a dataset of student responses that have been graded by human teachers. The model learns to identify key concepts such as "heat," "liquid," "vapor," and their relationships, which it uses to grade new responses.

Machine learning models find application in various contexts: In standardized testing scenarios, they are employed to assess open-ended questions, enhancing the grading process. Additionally, in formative assessments within educational settings, these models facilitate the provision of immediate feedback to students, aiding in gauging their comprehension of the material. Moreover, they are well-suited for large-scale assessments, particularly in environments like Massive Open Online Courses (MOOCs) and online learning platforms, where they can efficiently grade large volumes of responses.

### 4.2.1 Machine learning models have diverse applications:

In standardized testing, they excel at grading open-ended questions, enhancing the efficiency and accuracy of the assessment process. For formative assessments, these models enable instantaneous feedback to students, aiding in their comprehension and learning progression. Furthermore, they are particularly advantageous for large-scale assessments, effectively grading substantial volumes of responses in platforms such as MOOCs (Massive Open Online Courses) and online learning environments.

Machine learning models offer several benefits:

1. Scalability: They demonstrate the ability to efficiently manage large quantities of responses, facilitating streamlined grading processes.

2. Consistency: By mitigating human bias and variability, these models ensure uniform grading outcomes.

3. Adaptability: They possess the capability to refine and adjust to new grading criteria and patterns, enhancing their effectiveness over time.

However, they also present certain challenges:

1. Data Dependency: Effective functioning relies on extensive, high-quality datasets for training, posing a requirement that can be resource-intensive.

2. Black-Box Nature: The lack of transparency regarding decision-making processes can present challenges, making it difficult to understand how outcomes are determined.

3. Bias and Fairness: There is a risk of perpetuating biases inherent in the training data, potentially compromising the fairness and impartiality of grading outcomes.

### 4.2.2 Evolution of Machine Learning Models

Early machine learning models for automated grading used simpler algorithms and smaller datasets. Techniques such as Naive Bayes and decision trees were common, focusing on basic feature extraction like word frequency and simple syntactic patterns (Rudner et al., 2006).

An early ML model might use word frequency counts to grade responses to a question on historical events, recognizing key terms like "war," "treaty," and "independence."

As computational power and data availability increased, more advanced models were developed. These include:

1. Support Vector Machines (SVM): Used for classification tasks, SVM models can handle high-dimensional data and complex decision boundaries.

2. Neural Networks: Deep learning models, particularly neural networks, have shown significant improvements in understanding complex patterns and semantics in text.

For instance, a neural network model might analyze essays on climate change by understanding nuanced arguments, sentence structures, and the coherence of ideas presented.

Machine learning models have revolutionized automated grading by providing more accurate, scalable, and adaptable solutions compared to rule-based systems. They have enabled educational institutions to efficiently manage large-scale assessments and provide timely feedback to students, enhancing the overall learning experience (Attali & Burstein, 2006).

# 5. AI INTEGRATION: THE RISE OF GPT-4

The integration of artificial intelligence (AI) into automated grading systems has significantly advanced with the development of sophisticated models like GPT-4. These models leverage deep learning and natural language processing (NLP) to understand and grade complex student responses with unprecedented accuracy.

GPT-4 (Generative Pre-trained Transformer 4) is an advanced AI model developed by OpenAI. It represents a significant leap in natural language understanding and generation capabilities compared to its predecessors. GPT-4 is designed to comprehend and generate human-like text based on the input it receives, making it highly suitable for a range of applications, including automated grading.

GPT-4 boasts a range of key features that elevate its capabilities in educational assessment. Firstly, its language generation prowess enables the generation of detailed feedback for students, transcending the mere provision of grades to offer valuable insights for learning improvement. Furthermore, concerted efforts have been directed towards bias mitigation within GPT-4, resulting in assessments that strive for fairness and impartiality.

Underpinning its functionality is a sophisticated deep learning architecture, specifically a transformer architecture, renowned for its effectiveness in tasks involving sequence-to-sequence operations like text generation and comprehension. Coupled with this architecture is the extensive large-scale training undergone by GPT-4, which equips it to comprehend and generate text across a multitude of topics and contexts, ensuring versatility and adaptability in its application.

Central to GPT-4's efficacy is its contextual understanding capability, enabling it to grasp the context of student responses with precision. This contextual comprehension extends to nuances and implicit meanings within text, rendering GPT-4 adept at handling diverse and complex student submissions. Moreover, its adaptability shines through in its capacity for fine-tuning tailored to specific tasks, including automated grading, thus enhancing both its accuracy and relevance in educational assessment scenarios.

## 5.1 Applications in Automated Grading

GPT-4's advanced capabilities have made it a powerful tool for automated grading systems. Its applications span across various types of assessments, providing detailed feedback and accurate grading.

GPT-4 can evaluate short answer questions by understanding the context and content of the responses, comparing them to expected answers, and assigning appropriate grades. For instance, a question like "Explain the significance of the water cycle," GPT-4 can assess the student's response by identifying key elements such as evaporation, condensation, and precipitation, and understanding the explanation's coherence and completeness.

GPT-4 excels can also grade essays by analyzing structure, argumentation, grammar, and adherence to the prompt. It can provide holistic scores as well as detailed feedback on different aspects of the writing. For example, an essay about the impacts of global warming, GPT-4 can evaluate the clarity of the argument, the use of evidence, the organization of ideas, and the quality of writing, offering both a grade and constructive feedback.

GPT-4 can be integrated into educational platforms to provide real-time feedback to students as they write. This immediate feedback helps students improve their responses before final submission. This can be illustrated with a student writes an essay on renewable energy, GPT-4 can suggest improvements in sentence structure, point out missing arguments, and recommend additional evidence, enhancing the learning experience.

In the realm of automated grading, GPT-4 presents several distinct advantages. Foremost among these is its remarkable accuracy, stemming from its profound comprehension of language. Unlike traditional rule-based systems or simpler machine learning models, GPT-4 can evaluate responses with a higher degree of precision, leading to more reliable grading outcomes.

Moreover, GPT-4 offers consistency in grading, a notable departure from the variability often encountered with human graders. By providing uniform assessments across a vast number of responses, it ensures fairness and impartiality in the evaluation process.

Scalability is another area where GPT-4 shines. Its ability to efficiently handle large volumes of responses makes it particularly well-suited for large-scale assessments, including standardized tests and Massive Open Online Courses (MOOCs). This scalability not only enhances efficiency but also enables educators to manage assessments on a broader scale effectively.

One of the most significant advantages of GPT-4 in automated grading is its capacity to offer detailed feedback. Beyond assigning grades, GPT-4 can provide nuanced and constructive feedback to students. This feedback aids in understanding mistakes and areas for improvement, ultimately fostering a more conducive learning environment.

While GPT-4 offers numerous advantages in automated grading, its integration into grading systems also poses significant challenges:

One notable challenge is the potential for bias. Despite efforts to mitigate biases, GPT-4 may inadvertently perpetuate biases present in the training data, resulting in unfair grading outcomes that disadvantage certain groups.

Another issue is the interpretability of GPT-4's decision-making process. As a complex deep learning model, it operates with a level of opacity, making it challenging to discern how specific grades are determined. This lack of transparency can undermine trust in the grading system and raise concerns about accountability.

Furthermore, the accuracy and effectiveness of GPT-4 are heavily reliant on the quality and diversity of the training data. If the training data is not sufficiently comprehensive or representative, GPT-4 may struggle to accurately assess responses, leading to inconsistencies and inaccuracies in grading.

Addressing these limitations and challenges is crucial for ensuring the ethical and reliable use of GPT-4 in automated grading systems, requiring ongoing efforts to enhance fairness, transparency, and data quality.

The emergence of GPT-4 has ushered in a transformative era in educational assessment, leaving a profound impact on various facets of education:

One significant effect is the enhancement of learning outcomes. GPT-4's ability to offer detailed and prompt feedback

empowers students to learn more effectively and make substantial improvements in their academic performance. This personalized feedback aids in addressing specific areas of weakness, promoting deeper understanding, and fostering continuous learning.

Moreover, GPT-4 contributes to increased efficiency in the educational landscape. By automating the grading process, it alleviates the grading workload for educators, enabling them to allocate more time and resources to teaching and fostering meaningful interactions with students. This shift allows educators to focus on delivering high-quality instruction and providing individualized support to students, ultimately enriching the educational experience.

Additionally, the integration of GPT-4 in educational assessment holds the potential to advance equity in education. When appropriately calibrated and monitored, GPT-4 can help standardize grading practices and mitigate the influence of human bias. By offering impartial evaluations, it promotes fairness and equal opportunities for all students, regardless of background or demographic factors. This fosters a more inclusive learning environment where every student has the opportunity to succeed based on merit.

Overall, the advent of GPT-4 marks a significant step forward in educational assessment, promising to revolutionize teaching and learning practices while promoting fairness, efficiency, and equity in education.

## 5.2 Implementation in Grading Systems

GPT-4 can be integrated into grading systems to automate the evaluation of short answer responses. This integration involves several steps:

1. Data Preprocessing: Responses are preprocessed to ensure they are in a format suitable for analysis by GPT-4.

2. Model Training: GPT-4 is fine-tuned using a dataset of graded responses to improve its accuracy in specific educational contexts.

3. Evaluation and Feedback: The model grades responses and generates feedback, which can be reviewed by educators for accuracy.

## 5.3 Challenges and Considerations

While AI-driven grading systems offer many benefits, several challenges remain. These include ensuring the privacy and security of student data is paramount in the implementation of GPT-4 for educational assessment. Protecting sensitive information from breaches and unauthorized access is critical to maintaining student trust and complying with legal and ethical standards.

Bias mitigation remains an ongoing necessity. Continuous efforts are required to identify, address, and reduce biases inherent in AI models to ensure fair and equitable grading. This involves rigorous testing, validation, and refinement of AI systems to prevent the perpetuation of existing disparities.

Gaining acceptance and trust from educators and students is crucial for the effectiveness of AI-driven grading systems. Both groups need to be confident in the accuracy, fairness, and reliability of these systems. Building this trust involves transparency in how the AI operates, demonstrating its benefits, and providing support for its integration into educational practices.

## 6. CONCLUSION

The evolution of short answer grading systems has been marked by significant technological advancements, from manual grading methods to the sophisticated AI-driven solutions of today. GPT-4 represents a major milestone in this journey, offering highly accurate, efficient, and fair grading capabilities. As technology continues to evolve, the future of short answer grading systems holds great promise, with the potential to further enhance educational assessments and improve learning outcomes for students.

## 7. REFERENCES

[1] M. A. Hearst, "The debate on automated essay grading," IEEE Intelligent Systems and their Applications, vol. 15, no. 5, pp. 22-37, 2000.

[2] E. B. Page, "The use of the computer in analyzing student essays," International Review of Education, vol. 14, no. 2, pp. 210-225, 1968.

[3] S. Valenti, F. Neri and A. Cucchiarelli, "An overview of current research on automated essay grading.," Journal of Information Technology Education: Research,, vol. 2, pp. 319-332, 2003.

[4] M. D. Shermis and J. Burstein, "Automated essay scoring: A cross-disciplinary perspective," Lawrence Erlbaum Associates, 2003.

[5] L. M. Rudner, V. Garcia and C. Welch, "An evaluation of IBM's SPSS Modeler automated essay scoring system.," The Journal of Technology, Learning, and Assessment, vol. 4, no. 7, 2006.

[6] Y. Attali and J. Burstein, "Automated essay scoring with e-rater," The Journal of Technology, Learning and Assessment, vol. 4, no. 3, 2006.

[7] L. Anderson and D. Krathwohl, A Taxonomy for Learning, Teaching, and Assessing: A Revision of Bloom's Taxonomy of Educational Objectives, Longman, 2001.

[8] D. Sadler, "Formative assessment and the design of instructional systems," Instructional Science, vol. 18, no. 2, pp. 119-144, 1989.

[9] N. Gronlund, Assessment of Student Achievement, Allyn & Bacon, 2006.

[10] J. Wang and M. S. Brown, "Automated essay scoring versus human scoring: A comparative study," Journal of Educational Computing Research, vol. 35, no. 3, pp. 325-348, 2007.

[11] M. D. Shermis and J. Burstein, Automated Essay Scoring: A Cross-disciplinary Perspective, Lawrence Erlbaum Associates, 2003.

[12] S. Burrows, I. Gurevych and B. Stein, "The Eras and Trends of Automatic Short Answer Grading," *International Journal of Artificial Intelligence in Education, vol. 25, no. 1, pp. 60-117, 2015.

[13] S. Jordan and T. Mitchell, "e-Assessment for learning? The potential of short-answer free-text questions with tailored feedback," British Journal of Educational Technology, vol. 40, no. 2, pp. 371-385, 2009.

[14] A. e. a. Radford, "Language Models are Unsupervised Multitask Learners," OpenAI, 2019.

[15] T. e. a. Brown, Language Models are Few-Shot Learners, arXiv preprint arXiv:2005.14165, 2020.

[16] OpenAI, "GPT-4 Technical Report," OpenAI, 2023.

[17] R. Bennett and M. Zhang, "Validity and Reliability in Automated Essay Scoring.," Educational Measurement: Issues and Practice, vol. 35, no. 1, pp. 2-12, 2016.

[18] E. e. a. Higgins, "The future of AI in Education: Insights from the NMC Horizon Report," International Journal of Artificial Intelligence in Education, vol. 29, no. 2, pp. 123-142, 2019.