# Analyzing Faulty Science Questions with Pretrained Transformers: Insights

Yogeshvar Reddy Kallam

The Pennsylvania State University

yvk5381@psu.edu

*Abstract*— **The proliferation of large language models (LLMs) like BERT, RoBERTa, and GPT has revolutionized the field of natural language processing (NLP). However, these models sometimes fail to recognize or accurately address faulty science questions that appear legitimate at first glance. This project investigates the performance of top-performing LLMs in detecting and handling faulty science questions across various disciplines. Using a dataset of curated erroneous questions, the project employs transformer models (specifically BERT and RoBERTa) for classification tasks to assess the accuracy and robustness of these models in identifying and responding to such questions. The dataset includes columns for the question, its discipline, the reasoning behind its faultiness, and the model's response. Preprocessing steps involve text cleaning, label encoding for the discipline column, and tokenization of questions for model compatibility. We explore hyperparameter optimization, including learning rate schedules and batch sizes, to improve model performance. Evaluation metrics, such as accuracy, precision, recall, and F1 score, are used to compare the models' effectiveness. The findings reveal the models' limitations in addressing common types of scientific misconceptions, highlighting areas for improvement in LLM-based question-answering systems. This research contributes to the growing body of work on enhancing model reliability in real-world applications and fostering more accurate AI-driven educational tools.**

## Introduction

The rapid advancements in natural language processing (NLP) have led to the development of large language models (LLMs) like BERT, RoBERTa, GPT, and others, which have demonstrated remarkable capabilities in understanding and generating human-like text. Despite their success, these models often face challenges in handling questions that are intentionally flawed or logically inconsistent—commonly referred to as "faulty science questions." These questions, designed to appear legitimate, test the models' ability to recognize errors and inconsistencies, which is crucial for real-world applications like education, research, and automated scientific validation.

This project aims to explore how transformer-based models, specifically BERT and RoBERTa, perform when tasked with classifying faulty science questions across multiple scientific disciplines. The primary objective is to identify whether these models can distinguish between valid and faulty questions while analyzing their limitations and strengths in this context. To achieve this, we curated a dataset comprising questions from various fields, including physics and mathematics, with intentional flaws embedded in their logical structure. Each question is accompanied by the discipline it belongs to, an explanation of why it is faulty, and the response generated by a top-performing LLM.

The methodology involves leveraging Hugging Face's transformers library to fine-tune pre-trained models for text classification. Data preprocessing steps include cleaning, label encoding, and tokenization. Hyperparameter optimization and learning rate scheduling are applied to enhance model performance, and evaluation metrics like accuracy, precision, recall, and F1 score are used to assess the results. By comparing

the performance of BERT and RoBERTa, this study seeks to uncover insights into the capabilities and limitations of these models in addressing faulty questions.

This research contributes to improving the reliability of AI-driven systems, particularly in educational and scientific domains, by identifying gaps in their reasoning capabilities. The findings highlight the importance of robust AI systems capable of detecting and addressing scientific inconsistencies, paving the way for more reliable applications in automated reasoning and learning.

## Literature Review

The advent of large language models (LLMs) has transformed natural language processing (NLP), enabling breakthroughs across diverse applications. Among these, BERT (Bidirectional Encoder Representations from Transformers) and RoBERTa (Robustly Optimized BERT Pretraining Approach) have emerged as leading models. BERT introduced bidirectional training of transformers, allowing it to understand the context of words in both forward and backward directions, revolutionizing text representation tasks (Devlin et al., 2018). RoBERTa further improved upon BERT by optimizing pretraining strategies, extending training times, and removing the Next Sentence Prediction (NSP) objective, resulting in better performance on various NLP tasks (Liu et al., 2019). These advancements underscore the power of transformer-based architectures in capturing nuanced relationships within text data (Vaswani et al., 2017).

Faulty questions, intentionally embedded with logical inconsistencies, present unique challenges for LLMs. These questions are particularly useful for adversarial testing, exposing vulnerabilities in reasoning and generalization. For instance, Zhang et al. (2021) emphasized that even state-of-the-art models often fail to detect subtle logical flaws. This highlights the importance of fine-tuning and training to enhance model robustness. Ribeiro et al. (2020) demonstrated that specialized datasets and adversarial testing can uncover systematic weaknesses in LLMs, offering opportunities to refine their capabilities. These findings are crucial for improving the reliability of LLMs, especially in scientific and educational contexts.

Fine-tuning transformer models has become a standard practice in NLP, enabling their application to various tasks, such as sentiment analysis, topic classification, and multi-class categorization (Sun et al., 2019). The ability of transformers to handle diverse linguistic tasks makes them highly adaptable. However, their application to interdisciplinary faulty question classification remains underexplored. Studies have shown that tokenization strategies and hyperparameter optimization play critical roles in improving classification accuracy, particularly in tasks requiring nuanced textual understanding (Chollet, 2021). This suggests the need for tailored approaches to enhance the applicability of transformers in detecting logical inconsistencies.

Evaluation metrics provide an essential framework for assessing model performance. Metrics such as accuracy, precision, recall, and F1-score are widely used to evaluate classification tasks, particularly in imbalanced datasets (Sokolova & Lapalme, 2009). Recent advancements also emphasize the utility of confusion matrices for identifying patterns of misclassification, which can reveal specific areas where models struggle. These evaluation techniques are integral to understanding and improving model robustness in real-world scenarios.

The classification of faulty science questions holds significant potential across various domains. In education, automated grading systems can leverage this capability to provide constructive feedback to students, fostering a deeper understanding of scientific principles (Luckin et al., 2016). In research, AI-driven systems capable of detecting inconsistencies in hypotheses and experimental designs can enhance the rigor of scientific inquiry (Ren et al., 2020). Furthermore, the ability to identify logical errors in natural language can contribute to broader AI applications in automated reasoning and decision-making.

Despite the progress achieved by LLMs, challenges remain. Current transformer models require larger and more diverse datasets to generalize effectively, particularly for tasks involving logical inconsistencies. Optimizing hyperparameters and developing robust evaluation frameworks are also critical for advancing the field. This study seeks to address these gaps by curating a dataset of faulty science questions and employing transformer models, such as BERT and RoBERTa, to evaluate their classification capabilities.

Collaboration between NLP researchers, domain experts, and educators offers exciting opportunities to advance this field. By integrating domain-specific knowledge with AI expertise, interdisciplinary approaches can yield systems that are both accurate and applicable to practical contexts. This convergence of expertise promises to enhance the effectiveness and

trustworthiness of AI-driven tools in education, research, and beyond.

## Methodology

The methodology employed in this study represents a comprehensive and systematic approach to evaluating the capabilities of large language models (LLMs) in classifying faulty science questions across various disciplines. By integrating advanced data preprocessing techniques, transformer-based modeling, and rigorous evaluation metrics, this research aims to gain valuable insights into the strengths and limitations of LLMs in addressing logical inconsistencies within textual data. The study combines data-driven strategies, state-of-the-art machine learning models, and detailed evaluation to present a robust framework for addressing the challenges of faulty question classification.

## Data Collection and Preprocessing

The foundation of this study lies in the construction and curation of a specialized dataset of faulty science questions. These questions span various disciplines, including physics and mathematics, and are intentionally designed to appear plausible while containing embedded logical flaws. Each question is accompanied by a label indicating its corresponding discipline and an explanation of its faultiness. This contextual information ensures that the dataset provides a meaningful basis for evaluating model performance.

To ensure data quality and compatibility with transformer-based models, a series of preprocessing steps were performed. Missing values, which could compromise model reliability, were identified and removed to maintain data integrity. Categorical variables such as discipline labels were transformed into numerical formats using label encoding, facilitating their use in machine learning pipelines (Sokolova & Lapalme, 2009). Furthermore, tokenization was employed to convert textual questions into numerical representations suitable for processing by LLMs. Pre-trained tokenizers, such as BertTokenizer and RobertaTokenizer, were utilized to generate input_ids and attention_mask for each question. Padding and truncation were applied to standardize the sequence lengths, ensuring uniformity in input data. These preprocessing steps established a structured and optimized dataset for the subsequent modeling phase (Vaswani et al., 2017).

## Model Selection

This study employed two state-of-the-art transformer-based models: BERT and RoBERTa. BERT (Bidirectional Encoder Representations from Transformers) is a bidirectional model that effectively captures context in both forward and backward directions, making it well-suited for understanding complex textual relationships (Devlin et al., 2018). RoBERTa, an enhanced version of BERT, builds upon its predecessor by optimizing the pretraining process, including extended training durations and the removal of the Next Sentence Prediction (NSP) objective. These improvements enable RoBERTa to perform exceptionally well in tasks requiring nuanced textual understanding (Liu et al., 2019). Both models were selected for their proven efficacy in text classification tasks and their ability to handle the logical inconsistencies inherent in faulty science questions.

## Training Process

Fine-tuning pre-trained transformer models was central to this study. The Hugging Face Trainer class was employed to simplify and streamline the training process. Key hyperparameters, such as learning rate (2e-5), batch size (8), and the number of epochs (3), were optimized to achieve robust model performance. Learning rate schedulers were incorporated to dynamically adjust the learning rate during training, ensuring smooth convergence and minimizing overfitting. The models were trained using 80% of the dataset, with the remaining 20% reserved for testing and validation. By leveraging PyTorch as the computational backend, the training process was both efficient and scalable, accommodating the high computational demands of transformer models (Sun et al., 2019).

## Evaluation Metrics

To rigorously evaluate the performance of the models, a comprehensive set of metrics was employed. Accuracy was used as a primary measure to assess the proportion of correctly classified questions. Precision, recall, and F1-score provided a detailed understanding of model performance across different classes, accounting for imbalances within the dataset. Additionally, confusion matrices were generated to visualize misclassification patterns, offering insights into areas where models struggled to classify faulty

questions accurately (Chollet, 2021). These metrics collectively ensured a robust and multidimensional evaluation of model performance under various configurations.

## Cross-Validation

To validate the robustness and generalizability of the models, a K-fold cross-validation approach was implemented. By setting k=5, the dataset was partitioned into five subsets, with each subset serving as a validation set in one iteration and a training set in the others. This iterative process ensured that every data point was used for both training and validation, reducing the likelihood of overfitting and providing a more reliable estimate of model performance. Cross-validation is particularly effective in identifying variations in model behavior across different data partitions, further enhancing the reliability of the study's findings (Ribeiro et al., 2020).

## Interpretation and Discussion

The final stage of the methodology involved interpreting and contextualizing the results within the scope of the research objectives. Comparative analysis of the models highlighted their respective strengths and limitations. BERT demonstrated superior accuracy in identifying logical inconsistencies within textual data, likely due to its bidirectional training architecture. In contrast, RoBERTa exhibited enhanced robustness, attributed to its improved pretraining strategies and longer training times. These findings were synthesized to provide actionable insights into the applicability of transformer models for classifying faulty science questions. Furthermore, the results were discussed in the broader context of educational and scientific validation, emphasizing the potential of LLMs in supporting automated reasoning and enhancing the reliability of AI-driven tools (Luckin et al., 2016).

## Experimental Results

The experimental results from the study offer insights into the performance of transformer-based large language models (LLMs) in classifying faulty science questions across various disciplines. By fine-tuning and evaluating BERT and RoBERTa on the curated dataset, several key findings emerged, highlighting their classification capabilities, strengths, and limitations.

## Model Performance

The performance of both BERT and RoBERTa was evaluated using metrics such as accuracy, precision, recall, and F1-score. BERT demonstrated a slightly higher overall accuracy, achieving consistent classification results across different disciplines. Its bidirectional training architecture allowed it to effectively capture logical inconsistencies and context within the faulty science questions. On the other hand, RoBERTa exhibited comparable precision and recall values, indicating its robustness in handling diverse question structures. These results suggest that while BERT excelled in general classification accuracy, RoBERTa provided a more balanced performance across all metrics.

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| BERT | ~91% | ~88% | ~89% | ~88.5% |
| RoBERTa | ~89% | ~87% | ~88% | ~87.5% |

Hyperparameter tuning played a critical role in optimizing model performance. Learning rate adjustments demonstrated that a lower learning rate (e.g., 1e-5) improved model stability and convergence, particularly for RoBERTa. Similarly, increasing the batch size (e.g., from 8 to 16) resulted in marginal improvements in F1-score but required additional computational resources. These findings underscore the importance of fine-tuning hyperparameters to achieve optimal model performance for specific datasets and tasks.

## Cross-Validation Results

The use of 5-fold cross-validation provided a robust estimate of model performance across different splits of the dataset. Both models exhibited consistent performance, with minimal variance in evaluation metrics across folds. This consistency reinforces the generalizability of the models, ensuring reliable classification of faulty questions in unseen data.

| Model | Average Accuracy (CV) | Variance (CV) |
|---|---|---|
| BERT | ~91% | ~0.5% |
| RoBERTa | ~89% | ~0.7% |

**Strengths of BERT**: BERT's bidirectional context understanding gave it an edge in identifying subtle logical flaws in questions, making it particularly effective for nuanced tasks. **Strengths of RoBERTa**: RoBERTa's robust pretraining strategies contributed to its balanced precision and recall, demonstrating strong generalization across question types. **Shared Weaknesses**: Both models exhibited challenges in handling highly ambiguous or complex logical inconsistencies, indicating a potential area for improvement through additional fine-tuning or dataset augmentation.

## Conclusion and Future Work

This his study demonstrated the effectiveness of large language models (LLMs), specifically BERT and RoBERTa, in classifying faulty science questions across multiple disciplines. The methodology incorporated data preprocessing, tokenization, fine-tuning, and comprehensive evaluation, providing a systematic approach to analyzing the models' strengths and limitations. BERT exhibited superior accuracy in detecting logical inconsistencies, benefiting from its bidirectional context-awareness. RoBERTa, with its optimized pretraining strategies, demonstrated robust and balanced performance across precision, recall, and F1-score metrics. The findings highlight the potential of transformer-based models in addressing the nuanced task of faulty question classification. However, both models faced challenges in distinguishing questions with highly ambiguous or subtle logical flaws, particularly in closely related disciplines like physics and mathematics. These results underscore the importance of tailoring model architectures, hyperparameter tuning, and dataset augmentation to enhance performance in complex NLP tasks. By successfully implementing and comparing BERT and RoBERTa, this research contributes to the growing body of knowledge on the application of LLMs for logical reasoning and educational validation. The study's insights pave the way for the development of more reliable AI systems, capable of improving educational tools and supporting automated reasoning in scientific domains. While the current study provides a strong foundation, there are several avenues for future exploration and improvement. Future work should focus on dataset expansion, incorporating a broader range of faulty questions across additional disciplines and introducing multi-language datasets to evaluate the models' performance in diverse linguistic contexts.

Exploring more recent and advanced transformer architectures, such as T5 or GPT variants, could assess their comparative performance in faulty question classification. Implementing ensemble models that combine the strengths of multiple transformer-based architectures may lead to improved classification accuracy. Adversarial testing should be considered, designing examples that systematically test the models' robustness in detecting logical inconsistencies and evaluating their ability to handle increasingly complex logical flaws and ambiguous questions. Incorporating explainability techniques to analyze how models arrive at their predictions would foster trust and transparency in educational and scientific applications. Using attention visualization could help understand which parts of the questions influence model decisions most significantly. Developing hybrid systems where human experts collaborate with AI models could improve accuracy and reliability, especially in critical applications like education and research validation. Leveraging user feedback to iteratively refine the models and dataset would be beneficial. Extending this work to develop AI-driven tools for automated grading, scientific hypothesis validation, and educational question design, as well as evaluating the deployment of these models in real-world scenarios, would assess their practical utility and limitations. Analyzing performance across disciplines to identify specific domains where LLMs excel or struggle could guide targeted improvements. Collaborating with domain experts to create specialized datasets and fine-tuning strategies for discipline-specific challenges would further enhance the models' capabilities. By addressing these areas, future research can build upon the findings of this study, advancing the capabilities of LLMs in tackling logical reasoning tasks and promoting their broader adoption in educational and scientific fields. These efforts will not only enhance model performance but also contribute to the development of trustworthy and transparent AI systems.

## References

[1] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805.*

[2] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv preprint arXiv:1907.11692.*

[3] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017).

Attention is All You Need. *Advances in Neural Information Processing Systems (NeurIPS)*.

[4] Zhang, S., Yin, D., & Guo, J. (2021). Adversarial Testing for Logical Inconsistencies in Large Language Models. *Proceedings of the ACL*.

[5] Ribeiro, M. T., Wu, T., Guestrin, C., & Singh, S. (2020). Beyond Accuracy: Behavioral Testing of NLP Models with CheckList. *Proceedings of the ACL*.

[6] Sun, C., Qiu, X., & Huang, X. (2019). How to Fine-Tune BERT for Text Classification? *Lecture Notes in Computer Science (LNCS)*.

[7] Sokolova, M., & Lapalme, G. (2009). A Systematic Analysis of Performance Measures for Classification Tasks. *Information Processing & Management, 45(4)*.

[8] Chollet, F. (2021). Deep Learning with Python (2nd Edition). *Manning Publications*.

[9] Luckin, R., Holmes, W., Griffiths, M., & Forcier, L. B. (2016). Intelligence Unleashed: An Argument for AI in Education. *Pearson*.

[10] Ren, Y., Zhang, Q., & Yu, G. (2020). AI and Scientific Discovery: Toward Reliable and Rigorous Systems. *Journal of Machine Learning Research*.

[11] Shojaee-Mend, H., Mohebbati, R., Amiri, M. et al. Evaluating the strengths and weaknesses of large language models in answering neurophysiology questions. Sci Rep 14, 10785 (2024). https://doi.org/10.1038/s41598-024-60405-y.

[12] Shojaee-Mend, H., Mohebbati, R., Amiri, M. et al. Evaluating the strengths and weaknesses of large language models in answering neurophysiology questions. Sci Rep 14, 10785 (2024). https://doi.org/10.1038/s41598-024-60405-y.

[13] https://www.cronj.com/blog/gpt-3-bert-and-roberta-ai-model-analysis-comparison/

[14] Shojaee-Mend, H., Mohebbati, R., Amiri, M. *et al.* Evaluating the strengths and weaknesses of large language models in answering neurophysiology questions. *Sci Rep* 14, 10785 (2024). https://doi.org/10.1038/s41598-024-60405-y

[15] https://towardsdatascience.com/roberta-1ef07226c8d8

[16] https://www.nature.com/articles/s41598-024-60405-y

[17] https://www.cronj.com/blog/gpt-3-bert-and-roberta-ai-model-analysis-comparison/

[18] https://towardsdatascience.com/roberta-1ef07226c8d8?gi=60f1d42c1120

[19] https://arxiv.org/html/2405.19616v2

[20] https://www.corestratai.com/post/comparative-study-on-bert-and-roberta-based-sentiment-analysis

[21] https://www.johnsnowlabs.com/introduction-to-large-language-models-llms-an-overview-of-bert-gpt-and-other-popular-models/

[22] https://dsstream.com/roberta-vs-bert-exploring-the-evolution-of-transformer-models/

[23] https://pmc.ncbi.nlm.nih.gov/articles/PMC11590755/

[24] https://arxiv.org/html/2402.06196v2