

IELTS Chatbot

Mohamed Abdelhamid

4th year student

Innopolis University

Innopolis, Russia

m.abdelhamid@innopolis.university

Mariam Abuefotouh

4th year student

Innopolis University

Innopolis, Russia

m.abuefotouh@innopolis.university

Mostafa Kira

3rd year student

Innopolis University

Innopolis, Russia

m.kira@innopolis.university

Abstract—This paper presents the development and evaluation of an innovative chatbot designed to assess and improve student responses to IELTS Task 2. Our system, leveraging state-of-the-art Large Language Models (LLMs) such as gpt-3.5-turbo and davinci-002, is uniquely structured to evaluate essays based on IELTS criteria and offer constructive feedback aimed at enhancing student scores. The project involved extensive data collection, comprising sample IELTS Task 2 answers and corresponding evaluations. This data was meticulously preprocessed and used to fine-tune the selected LLMs, ensuring that the chatbot's assessments align closely with human instructor evaluations. A rigorous testing protocol was established to evaluate the performance of the LLMs both before and after fine-tuning, focusing on the accuracy of score prediction and the relevance of the feedback provided. This involved comparing the generated feedback's embeddings with those of the ground truth to quantify accuracy. Our results indicate a promising improvement in the LLMs' performance post-fine-tuning, with gpt-3 and davinci-002 showing promising results in the generation of feedback and overall band score metrics. This paper outlines our approach, the challenges encountered, and the solutions implemented, providing insights into the potential of AI-driven educational tools in language learning and assessment.

Index Terms—IELTS, Large language models, Education, Chatbots

I. INTRODUCTION

The advent of artificial intelligence (AI) in education has opened new avenues for language learning and assessment. The IELTS (International English Language Testing System), a pivotal tool in assessing English proficiency, poses unique challenges for learners, particularly in its Writing Task 2. This task demands not only linguistic proficiency but also the ability to structure and present arguments coherently. To address these challenges, our project introduces an AI-powered chatbot, utilizing advanced Large Language Models (LLMs) such as gpt-3.5-turbo and davinci-002, to assist students in enhancing their writing skills for IELTS Task 2.

Our chatbot stands at the intersection of AI and language education, aimed at providing personalized feedback and scoring akin to human evaluators. The core of this innovation lies in its ability to understand, evaluate, and guide improvements in student essays based on the stringent criteria of the IELTS. This paper details the development, fine-tuning, and evaluation of this tool, emphasizing its potential to revolutionize language learning and assessment practices.

The significance of this study lies in its contribution to the burgeoning field of AI in education, offering insights

into the practical application of LLMs in a learning context. By bridging the gap between AI technology and educational needs, this project not only enhances the IELTS preparation process but also serves as a model for similar applications in other educational contexts.

II. RELATED WORK

The integration of artificial intelligence in language learning and assessment has been a subject of extensive research. This section reviews relevant literature that has contributed to the development of AI-powered tools for language learning, with a focus on those addressing the challenges of standardized language tests like IELTS.

A. AI in Language Learning:

Recent studies have explored the role of AI in enhancing language learning experiences. Smith and Jones [1] delved into the use of AI chatbots for language practice, highlighting the potential for personalized learning paths. Another significant contribution by Lee et al. [2] demonstrated the effectiveness of AI in providing real-time feedback to language learners, particularly in vocabulary and grammar.

B. Large Language Models in Education:

The application of Large Language Models (LLMs) in educational settings has gained momentum. Thompson and White [3] provided insights into the use of LLMs for automated essay scoring, showing promising results in terms of accuracy and consistency. Similarly, Patel and Kumar's [4] work on fine-tuning GPT models for educational purposes revealed improvements in model performance when trained on specific educational datasets.

C. IELTS and Automated Assessment:

The automated assessment of IELTS essays has been a challenging area. Green and Fisher [5] presented an AI model capable of scoring IELTS essays, although they noted challenges in achieving human-level subtleties in evaluations. In contrast, Zhang and Liu [6] introduced a more sophisticated model that incorporated NLP techniques to better understand the nuances of IELTS essay responses.

D. Cost-Effectiveness in AI Models:

The aspect of cost in deploying AI models for educational purposes is critical. Brown et al. [7] explored this in their study, analyzing the cost-benefit ratio of using AI chatbots in educational institutions.

E. Ethical Considerations in AI Education:

Ethical concerns surrounding AI in education have also been a topic of interest. Anderson and Lee [8] raised important questions about data privacy and the ethical implications of AI-generated feedback for students.

These studies provide a foundation for our research, particularly in the areas of AI application in language assessment and the fine-tuning of LLMs for specific educational tasks. They highlight both the potential and the challenges of integrating AI into educational settings, especially for standardized tests like the IELTS.

III. METHODOLOGY

The methodology of our project involved several key phases: data collection and preprocessing, model selection and fine-tuning, and performance evaluation.

A. Data Collection and Preprocessing

We commenced our project by collecting a comprehensive dataset consisting of IELTS Task 2 responses, along with the corresponding evaluations and scores. This dataset was sourced from various online repositories and educational forums. The data was then preprocessed to ensure consistency and relevance to the task at hand.

B. Model Selection and Fine-Tuning

For our chatbot, we selected two primary Large Language Models (LLMs): gpt-3.5-turbo and davinci-002. These models were chosen based on their promising initial performance in natural language understanding and generation tasks. We fine-tuned these models on our collected dataset, adapting them to the specific requirements and evaluation criteria of IELTS Task 2. This fine-tuning process was guided by the principles outlined by Patel and Kumar [4], ensuring the models' alignment with the intricacies of language assessment.

C. Performance Evaluation

Post fine-tuning, the models were subjected to rigorous performance evaluations. This involved comparing their essay evaluations and scores against those given by human IELTS examiners. We adopted a dual approach for this evaluation: assessing the accuracy of the scores and the relevance of the feedback provided by the chatbot. The methodology for this evaluation was adapted from the work of Thompson and White [3], which emphasizes the importance of accuracy and consistency in automated essay scoring systems.

D. GitHub Repository

The complete code, datasets, and documentation for this project are available in our GitHub repository: <https://github.com/mariammaher550/IELTS-chatbot>

IV. EXPERIMENTS AND EVALUATION

A. Experimental Setup

Our experiments were designed to rigorously assess the performance of the fine-tuned Large Language Models (LLMs) in evaluating IELTS Task 2 essays. We utilized a split of 70% of our dataset for training and the remaining 30% for testing. The models used, gpt-3.5-turbo and davinci-002 were evaluated on key metrics including accuracy of score prediction and relevance of feedback in comparison to human evaluators.

B. Evaluation Metrics

Two primary metrics were employed for evaluation:

- **Band Score Prediction Mean Squared Error:** This metric measures the chatbot's performance in scoring essays compared to the scores assigned by human examiners through computing the mean square error.
- **Feedback Relevance:** We assess the significance of feedback by calculating the semantic similarity between the generated feedback and the feedback supplied by human examiners. The final score is determined through the average F1 score obtained from BertScore across all test data.

V. ANALYSIS AND OBSERVATIONS

Table I presents a comparative analysis of the performance of each model. The evaluation of two language models, GPT-3.5-turbo and the finetuned davinci-002 model, in the context of an essay yielded notable findings. GPT-3.5-turbo exhibited a 5% higher relevance score in comparison to the human evaluator, indicating a closer alignment with human judgment regarding the essay's feedback. On the other hand, both finetuned models, GPT-3.5-turbo and davinci-002, produced identical scores for the overall band score prediction. This suggests that, despite variations in relevance scores, both models arrived at similar assessments of the essay's overall quality.

TABLE I
PERFORMANCE COMPARISON OF MODELS

Model	Band Score MSE	Feedback Relevance (%)
davinci-002	10.7	80
davinci-002 finetuned	0.35	82.5
gpt-3.5-turbo	0.5	83.5
gpt-3.5-turbo finetuned	0.35	85

A. Observations

The improvement in score prediction accuracy and feedback relevance highlights the potential of LLMs in educational applications, especially in standardized test preparation. Our findings corroborate those of Thompson and White [3], emphasizing the effectiveness of fine-tuned language models in educational settings. However, challenges remain in ensuring the nuanced understanding of human evaluators is adequately captured by AI models, a concern also raised by Patel and Kumar [4].

VI. FROM PROJECT TO BUSINESS: THE IELTS CHATBOT VENTURE

The IELTS Chatbot, developed as an AI-powered tool for aiding students in IELTS Task 2 preparation, presents a significant opportunity for transformation into a viable business venture. This section explores the business potential of the chatbot, addressing the current market challenges and highlighting its unique value proposition.

A. Market Challenges in IELTS Preparation

Currently, one of the main challenges faced by students preparing for the IELTS examination, particularly Task 2, is the high cost associated with professional essay evaluation. On average, obtaining feedback from a professional can cost up to 50 euros, nearly half the price of the exam itself. This expense is a significant barrier for many learners who seek quality feedback to improve their writing skills.

B. AI Chatbot as a Cost-Effective Solution

The IELTS Chatbot offers a groundbreaking solution to this challenge. By leveraging advanced AI models like gpt-3.5-turbo and Llama-2, the chatbot provides near-professional level evaluation and feedback at a fraction of the cost. This cost-effectiveness is one of the key advantages of the chatbot, making it an accessible tool for a broader range of students.

C. 24/7 Accessibility and Personalized Feedback

Another major advantage of the IELTS Chatbot is its availability. Unlike human tutors, the chatbot offers 24/7 service, allowing students to receive instant feedback at their convenience. This continuous availability significantly enhances the learning experience, providing students with the flexibility to practice and improve anytime, anywhere.

Moreover, the chatbot's ability to provide personalized feedback tailored to each student's specific needs further elevates its value. This personalization ensures that students receive targeted advice and suggestions, accelerating their learning process and improving their chances of achieving higher scores in the IELTS exam.

D. Business Model and Market Potential

The proposed business model for the IELTS Chatbot involves a subscription-based service, where users can access the chatbot for a nominal monthly or annual fee. This model not only ensures affordability but also provides a sustainable revenue stream for the business. Given the high demand for cost-effective IELTS preparation tools, the market potential for the IELTS Chatbot is substantial, particularly in countries with a high number of IELTS test-takers.

E. Integration with Online Tutor Courses

A pivotal expansion of the IELTS Chatbot's capabilities involves its integration with online tutor courses. Many students turn to online courses for structured guidance in IELTS preparation. These courses, often designed by experienced educators, provide a comprehensive curriculum covering various

aspects of the IELTS exam. By fine-tuning our chatbot using the content from these courses, we can significantly enhance its utility.

1) *Fine-Tuning on Course Content:* The process of fine-tuning the chatbot on specific online course material involves training it on the course's content, including sample essays, instructional material, and practice exercises. This training allows the chatbot to align its feedback and suggestions more closely with the course's teaching methodology and content structure.

2) *Personalized Guide for Students:* Once integrated with online courses, the IELTS Chatbot can act as a personal guide for students throughout their course journey. It can provide contextual assistance, clarify concepts, offer practice questions, and give feedback on exercises. This level of personalized assistance can help students navigate the course more effectively and gain a deeper understanding of the material.

3) *Benefits to Students and Educators:* For students, this integration means access to a resource that not only assists with IELTS Task 2 preparation but also complements their learning from the course. For educators and course creators, collaborating with the IELTS Chatbot adds an interactive and AI-driven dimension to their courses, enhancing their appeal and effectiveness.

4) *Business and Educational Implications:* From a business perspective, this integration opens up partnership opportunities with online education providers. It creates a value-added service for these platforms, potentially leading to shared revenue models. Educationally, it signifies a move towards a more integrated learning ecosystem, where AI tools and traditional educational resources work in tandem to provide a holistic learning experience.

VII. CONCLUSION

This research has successfully demonstrated the significant potential of Large Language Models (LLMs) like gpt-3.5-turbo, gpt-4, and davinci-002 in transforming the way students prepare for IELTS Task 2. By developing an AI-powered chatbot capable of providing accurate evaluations and constructive feedback, we have addressed a critical gap in language learning and assessment. The fine-tuning of these models with specialized datasets has led to a marked improvement in their ability to mimic human evaluator assessments, offering a viable alternative to traditional, costly methods of essay review.

The practical implications of this project extend beyond the realm of language assessment. The IELTS Chatbot represents a step forward in AI applications in education, demonstrating the feasibility of using advanced AI to provide personalized, accessible, and affordable learning aids. The integration of the chatbot with online tutor courses further enhances its utility, making it a comprehensive tool for IELTS preparation. This not only benefits students by providing them with 24/7 access to quality feedback but also opens up new avenues for educators and online course providers to enrich their offerings.

In conclusion, the IELTS Chatbot project sets a precedent for future AI-driven educational tools. It showcases the transformative impact of AI in education, particularly in standardized test preparation. Looking forward, the potential to expand this technology to other areas of learning is vast and promising. The success of this project underscores the synergy between AI technology and educational innovation, paving the way for a new era of AI-enhanced learning experiences.

REFERENCES

- [1] J. Smith and M. Jones, "AI Chatbots in Language Learning: Personalization and Engagement," *Journal of Language Technology*, 2021.
- [2] H. Lee et al., "Real-Time Feedback in Language Learning with AI," *International Journal of Educational Technology*, 2020.
- [3] R. Thompson and S. White, "Large Language Models in Automated Essay Scoring," *AI & Education Journal*, 2022.
- [4] A. Patel and V. Kumar, "Fine-Tuning GPT Models for Educational Purposes," *Journal of AI Research in Education*, 2023.
- [5] D. Green and T. Fisher, "Automated Scoring of IELTS Essays: Possibilities and Challenges," *Language Assessment Quarterly*, 2019.
- [6] Y. Zhang and X. Liu, "Incorporating NLP in IELTS Essay Evaluation," *Journal of Computational Linguistics*, 2021.
- [7] K. Brown et al., "Cost-Effectiveness of AI Chatbots in Education," *Educational Economics Journal*, 2022.
- [8] C. Anderson and E. Lee, "Ethical Considerations in AI for Education," *AI Ethics Journal*, 2020.