# GenAI Content Detection Task 2:
# AI vs. Human – Academic Essay Authenticity Challenge

**Shammur Absar Chowdhury**[1]**, Hind Almerekhi**[1]**, Mucahid Kutlu**[2]**, Kaan Efe Keleş**[3]**,**
**Fatema Ahmad**[1]**, Tasnim Mohiuddin**[1]**, George Mikros**[4]**, Firoj Alam**[1]

[1]Qatar Computing Research Institute, HBKU, Qatar, [2]Qatar University, Qatar
[3]TOBB ETU, Türkiye, [4]Hamad Bin Khalifa University, Qatar
{shchowdhury, halmerekhi, fakter, mmohiuddin, gmikros, fialam}@hbku.edu.qa
mucahidkutlu@qu.edu.qa, kaanefekeles@etu.edu.tr

## Abstract

This paper presents a comprehensive overview of the first edition of the *Academic Essay Authenticity Challenge*, organized as part of the GenAI Content Detection shared tasks collocated with COLING 2025. This challenge focuses on detecting machine-generated *vs* human-authored essays for academic purposes. The task is defined as follows: *"Given an essay, identify whether it is generated by a machine or authored by a human."* The challenge involves two languages: English and Arabic. During the evaluation phase, 25 teams submitted systems for English and 21 teams for Arabic, reflecting substantial interest in the task. Finally, seven teams submitted system description papers. The majority of submissions utilized fine-tuned transformer-based models, with one team employing Large Language Models (LLMs) such as Llama 2 and Llama 3. This paper outlines the task formulation, details the dataset construction process, and explains the evaluation framework. Additionally, we present a summary of the approaches adopted by participating teams. Nearly all submitted systems outperformed the n-gram-based baseline, with the top-performing systems achieving F1 scores exceeding 0.98 for both languages, indicating significant progress in the detection of machine-generated text.

## 1 Introduction

The rapid progress in Artificial Intelligence (AI) and the proliferation of generative content produced by LLMs have introduced transformative opportunities across various domains — yet they also pose profound challenges (Wu et al., 2023). One such challenge lies in the detection and prevention of misuse of LLMs in contexts such as fake news, misinformation, disinformation, and academic dishonesty (Tang et al., 2024). For instance, the volume of AI-generated news on misinformation-prone websites surged by 457% between January 1, 2022, and May 1, 2023, with a corresponding increase of 57.3% on mainstream platforms (Hanley and Durumeric, 2024). These issues pose substantial barriers to the broader adoption of LLMs, thereby limiting their potential across various applications. Effectively detecting LLM-generated content is crucial for leveraging the capabilities of these models while mitigating associated risks.

Researchers have responded to these challenges through a variety of approaches. Previous methods include classification algorithms designed to distinguish between AI-generated and human-authored text (Guo et al., 2023), as well as watermarking techniques (Szyller et al., 2021; He et al., 2022; Kirchenbauer et al., 2023). These watermarking approaches strategically embed imperceptible signatures within generated texts, enabling model-specific identification while maintaining human-indistinguishable quality. Other recent efforts have focused on the creation of question-answering datasets such as M4 (Wang et al., 2024b), generated by humans and ChatGPT in both English and Chinese and the associated shared task (Wang et al., 2024a).

Within academic settings, concerns surrounding the potential misuse of LLMs have intensified, particularly regarding academic dishonesty involving AI-assisted essay writing and problem-solving. Recent research has made considerable progress in the development of datasets and benchmarking efforts to address these issues. For instance, Yu et al. (2023) introduced the CHEAT dataset, which focuses on abstracts from IEEE Xplore, while Wang et al. (2024b) developed a comprehensive multilingual dataset. Additionally, Dugan et al. (2024) presented a robust dataset designed to address the challenge of detecting machine-generated text.

Despite these efforts, large-scale initiatives in academic contexts remain limited. Hence, this shared task aims to bridge this gap by tackling the task of distinguishing AI-generated essays from

human-authored ones. The challenge attracted substantial interest, with 99 teams registered to access the dataset and 56 teams actively participating in the development and evaluation phases. In the evaluation phase, 25 teams submitted systems for English, and 21 teams participated for Arabic. Furthermore, seven teams submitted system description papers. The majority of participating systems employed transformer-based models, while one team utilized state-of-the-art LLMs such as Llama 2 and Llama 3. Notably, most submissions outperformed the traditional n-gram-based baseline, signaling substantial progress in AI-generated content detection methodologies.

The subsequent sections of this paper are structured as follows: Section 2 provides a comprehensive review of related work. Section 3 presents the task formulation and dataset setup. Section 4 presents empirical results and offers a comprehensive overview of participating systems. Finally, Section 5 concludes with a summary of findings and future directions.

## 2  Related Work

The detection of AI-generated text relies on analyzing statistical patterns and linguistic features that distinguish human and machine writing styles. Zaitsu and Jin (2023) highlight that AI-generated text tends to use repetitive sentence patterns and a limited vocabulary, prioritizing clarity over the nuanced variations of human writing. Similarly, Weber-Wulff et al. (2023) report that such texts often exhibit lower syntactic complexity and reduced lexical diversity, making them identifiable through these markers. Additionally, Gallé et al. (2021) report that higher predictability in word $n$-gram is a key indicator of machine generated text.

Machine learning approaches have become central to AI-generated text detection. Darda et al. (2023) explored traditional classification algorithms such as Support Vector Machines (SVM) and Random Forest. Vora et al. (2023) propose a multimodal approach that uses BERT to analyze syntactic and semantic features of text and CNN architectures for image. Mikros et al. (2023) investigated using stylometric features and transformer-based models. Their findings showed that ensemble techniques, particularly those employing majority voting, outperformed individual classifiers.

There has also been effort to combine different machine learning approaches. For instance, deep learning architectures can extract features from text, while traditional classifiers make predictions based on these features, leveraging the strengths of both techniques (Bhattacharjee et al., 2023). Incorporating user feedback further enhances hybrid models, enabling them to adapt to real-world usage patterns (Rashidi et al., 2023).

Despite advancements in detection methodologies, significant limitations persist. Weber-Wulff et al. (2023) reveal that many detection tools struggle with high rates of false positives and false negatives, indicating a need for further refinement. According to Perkins et al. (2024), humans naturally incorporate varying sentence lengths and structures in their writing, creating what researchers call "burstiness"—a key feature that distinguishes human-authored content from AI-generated text. This variation in writing style, along with occasional grammatical inconsistencies and stylistic irregularities, represents the natural "imperfections" that make human writing unique. Interestingly, Liang et al. (2023) found that texts with lower levels of perplexity and coherence—characteristics often found in writing by non-native English speakers—are more likely to be flagged as human-authored.

Another challenge in AI-generated content detection is the lack of transparency in models' predictions, reducing their applicability in real-life scenarios, particularly in high-stakes contexts such as academia and forensic applications. Thus, a number of researchers worked on developing explainable AI (XAI) methods for AI generated text detection. For instance, Shah et al. (2023) develop an XAI model using stylistic features. Wu and Flanagan (2023) proposes a hybrid approach that combine statistical analysis with machine learning techniques. Additionally, the integration of user feedback into hybrid models may facilitate the development of more adaptive systems that can learn from usage patterns (Rashidi et al., 2023).

## 3  Task and Dataset

### 3.1  Task Definition

The main objective of the task is to detect whether the given candidate essay is AI-generated or human-written. Given the input essay $\mathbf{e}$, the task is to design a text detector $\mathcal{D}(\mathbf{e})$, such that the model outputs label indicating AI-generated or Human-authored content. For this edition, we designed the task as binary classification problem.

| | |
|---|---|
| **System Prompt** | You are a **{study_level}** student from **{country}**, preparing for the TOEFL exam. Your English proficiency level is **{proficiency_level}**. Your task is to write a well-structured TOEFL essay in response to the given prompt. Ensure your essay is clear and coherent, following the standard essay format: an introduction, body paragraphs, and a conclusion. Focus on presenting your ideas logically, using appropriate language, and providing relevant examples to support your arguments. Aim to demonstrate your proficiency in English through organized thought and effective communication. |
| **User Prompt** | Do you agree or disagree with the following statement: **"{statement}"** Write a well-structured essay expressing your opinion. Be sure to use specific reasons and examples to support your viewpoint. The essay should be between **{min_length}** and **{max_length}** words in length. Please provide only an essay and in a JSON object. No additional text or explanation. **{"essay": "your essay"}** |

Table 1: Example of *System* and *User Prompts* for training and validation in English essay generation. Similar prompts were used for Arabic essays. Variables include study_level ={'pre-university','university'}, proficiency_levels={'low','medium','high'}, country_list={'Arabic', 'German', 'French', 'Hindi', 'Italian', 'Japanese', 'Korean', 'Spanish', 'Telugu', 'Turkish', 'Chinese'}. For Arabic prompts, an additional variable, nativity={'native','non-native'} is used.

## 3.2 Datasets

The task aims to develop a system specifically designed for detecting AI generated text in academic essays. The dataset comprises essays authored by both native and non-native speakers, alongside AI-generated content. A significant challenge in this task was collecting authentic human-authored academic essays while addressing the following considerations:

- Ensuring author privacy, obtaining informed consent, and ethically sourcing the content.

- Verifying that the collected essays were genuinely authored by humans, free from any AI interference or plagiarism.

- Acquiring a diverse set of essays representing different academic levels and cultural backgrounds to ensure inclusivity in the dataset.

For the task, we focused on two languages: English and Arabic. For each language, we provided training, validation, dev-test, and the final test sets, which included human-authored and AI-generated texts. We released these data splits in two phases – *(i) Development phase* – we released the training, validation, and mock test data (dev-test); *(ii) Evaluation phase* – we released the final test set which

is used to rank the submitted system. Below, we discuss the dataset design for the development and final evaluation phases, respectively.

## 3.3 Development Phase

During the development phase we have released training, validation, and dev-test. For this phase, we first collected human-authored essays and essay topics. To create the data splits, we carefully designed each set to ensure unique essay topics, avoiding overlap between training, validation, and dev-test datasets.

Furthermore, within each split, we manually categorized the essay topics based on their thematic similarity. This classification is used to assign topics for generating essays using LLMs, and the rest is reserved exclusively for selecting human-authored essays from various existing datasets mentioned below. The final statistics of the dataset released in this phases are presented in Table 5.

**Human-authored Essay** The human-authored data was sourced from different language assessment datasets, including examinations like IELTS, and TOEFL among others. To ensure the authenticity of human-authored content, we selected essays that were either handwritten or composed in a supervised classroom setting, explicitly to make

sure that none of the texts were created with the assistance of generative technologies or online articles. This approach was designed to maintain the integrity of the datasets and accurately represent human academic writing.

For the English, we collected essay statements (essay prompt) and essays from:

- **IELTS Writing Scored Essays** Dataset[1] contains 1200 academic essays for varieties of prompts. Each essays are accompanied by the examiners' feedback along with scores

- **ETS Corpus of Non-Native Written English** corpus[2] contains 12,100 academic essays, written addressing eight different prompts, by non-native speakers from 11 different countries, as part TOEFL English proficiency exam. The dataset includes the speaker's native language along with scores they obtained for the corresponding essays. While the dataset was originally designed for native language identification tasks, its rich collection of academic essays, makes it highly suitable for supporting our AI-generated text detection efforts.

As for the Arabic subtask, the datasets we use are the following:

- **Arabic Learner Corpus (ALC)**[3] (Alfaifi and Atwell, 2013) includes 1,197 essays written by both native and non-native Arabic pre-university/university speakers from 67 nationalities. The dataset includes speakers' nationality along with the information if the essay was written in class or as homework. For the task, we only selected in-class essays, manually excluded off-topic essays, and reviewed the essays for any corrections.

- **Qatari Corpus of Argumentative Writing (QCAW) dataset**[4] (Zaghouani et al., 2024) is a collection of 195 argumentative essays written by native Arabic undergraduate students. The prompts given to the student were inspired by TOEFL writing exercises (Ahmed et al., 2023).

- **The CERCLL corpus**[5] includes $\approx 270$ essays written by non-native (L2) and heritage Arabic speakers.[6] The dataset includes information about the speakers' proficiency, along with the type – L2 *vs* heritage speakers. The dataset covers a wide range of topics and multiple genres, including description, narration, and instruction essays.

**AI-generated Essay** The generated essays, for both languages, utilized seven state-of-the art LLMs including: GPT-3.5-Turbo (2023-03-15-preview), GPT-4o (2024-08-06), GPT-4o-mini (2024-07-18) (OpenAI, 2024), Gemini-1.5 (Team, 2024), phi3.5,[7] Llama-3.1 (8B) (Abdin et al., 2024), and Claude-3.5.[8] To produce these essays, we designed the prompts by utilizing a selected subset of essay statements from the aforementioned datasets. The designed prompts included detailed instructions to emulate human writing styles, specify essay length requirements, and incorporate predefined personas reflecting various factors such as nativity and/or language proficiency, following the metadata and statistics obtained from the human-authored essay collections. This approach ensured the generation of essays that closely resemble real-world human writing in both style and content. An example of such a prompt is shown in Table 1.

### 3.4 Evaluation Phase

For the evaluation, we designed and developed a novel dataset, the **G**enerated and **R**eal **A**cademic **C**orpus for **E**valuation (GRACE), which includes both human-authored and AI-generated essays in English and Arabic.

### 3.4.1 Data Collection

For designing the human-authored portion of the dataset, we began by carefully designing test set essay statements aligned with those used in development phase topics. We selected five different essay types, and under each type, we created several essay statements (see Table 4 for examples). The topics include social influence & technology, lifestyle choices & preferences, cultural & global perspective, environmental & societal responsibility, and personal growth & experience.

---

```
You are tasked with generating creative and rigorous academic essays.
Here's how:
1) Topics Selection: You are provided with a set of topics: «<20 random topics»>. First, choose one
topic at random from this list.
2) Generate Related Topics: Based on the chosen topic, create 10 new topic ideas. These should be
different from the chosen topic but related in a way that someone interested in the initial topic
might also find these new ideas engaging.
3) Select Final Topic: From the 10 new topics, pick one at random to focus on.
4) Choose a Profession: List 10 random professions that are entirely unrelated to the final topic,
ensuring that they come from different fields or disciplines. These professions should be distinct
enough that their practitioners would not typically engage with or have knowledge about the topic.
Then, select one profession at random from this list.
5) Choose a Writing Style: List 10 distinct writing styles (e.g., persuasive, narrative, descriptive)
and choose one at random.
6) Essay Writing: Write an academic and creative essay on the chosen topic. This essay should be
written from the perspective of someone in the chosen profession and in the selected writing style.
Do not ever mention the chosen profession or writing style in the essay itself. Do not include
any personal opinions or experiences with regarding to the profession in the essay. Do not mention
anything about the chosen profession whatsoever.
Your output should be in JSON format, structured as follows:
{ "selected_topic": "<randomly selected topic from the given topics>", "generated_topics": [
"<generated topic 1>", "<generated topic 2>", "...", "<generated topic 10>" ], "final_topic":
"<randomly selected topic from generated_topics>", "professions": [ "<profession 1>", "<profession
2>", "...", "<profession 10>" ], "selected_profession": "<randomly selected profession
from professions>", "writing_styles": [ "<style 1>", "<style 2>", "...", "<style 10>" ],
"selected_writing_style": "<randomly selected style from writing_styles>", "essay": "<generated
essay>" }
Please proceed with this format to generate a fully structured JSON output. Remember to keep the
content diverse and creative throughout the process. The essay should be comprehensive, detailed,
and reflective of rigorous academic standards. The essay must be multiple paragraphs long (at least
1 page's worth). Return only the valid JSON output and nothing else. Good luck!
```

Table 2: **Freehand prompt** used to generate AI generated essays for the final test set.

**Essay Writing by Recruited Participants:** We then recruited[9] university students, both monolingual and bilingual, contribute to the essay writing. The participants were provided with a list of essay statements in their respective languages (either English or Arabic) and were asked to complete each essay within 30 minutes. They were instructed to limit the essays to 350–500 words and ensure they included an introduction, main arguments, and a conclusion. The essays must be written in Modern Standard Arabic (MSA) for Arabic, or in formal English for the English essays.

**Collected Essay Assignments:** Additionally, we collected previously submitted English *essay assignments* from university students to enrich the dataset.

*Anonymization of Personal Information* In the collected *essay assignments* we noticed that there were some information containing mentions of entities. Therefore, we anonymized them to ensure the removal of any information that could directly or indirectly identify the author or reveal any private infor-

mation about an entity that is not publicly known. This process was essential to uphold privacy standards and ethical considerations.

To achieve this, we followed these guidelines:

- *Author Identification Removal*: Any mention of names, addresses, affiliations, or specific details that could identify the essay's author was redacted.

- *Private Entity Information*: Any references to non-public entities, such as organizations, businesses, or private individuals mentioned in the essays, were removed or replaced with generic terms.

- *Sensitive Content*: Sensitive information, such as health conditions, financial details, or other personal data, was also removed to ensure privacy.

- *Consistency*: Replacement terms were standardized (e.g., "[NAME]", "[ADDRESS]", "[ORGANIZATION]") to maintain consistency throughout the dataset.

A team of five trained annotators was recruited

---

[9]We use a third-party company for the reward money. The amount was decided based on the standard local rate for data annotation.

```
Thoroughly rewrite the provided academic essay to enhance clarity, diversity in sentence structure,
and vocabulary richness, all while maintaining the original meaning and intent. Your goal is to
produce a refined and nuanced version of the text.
Aim to increase the essay's length by adding substantial elaborations, exploring various perspectives,
and providing comprehensive explanations that will offer a deeply layered and extensive output.
Deliver the output exclusively in JSON format with a single key "text" as shown below, ensuring that
no additional information or comments are included:
{{ "text": "<rewritten_and_greatly_expanded_academic_essay>" }}
Here is the passage to rewrite and extensively expand:
«<original_passage_start»> {the passage to be paraphrased} «<original_passage_end»>
```

Table 3: **Paraphrasing prompt** used to generate AI generated essays for the final test set.

| Question Type | Example Statements |
|---|---|
| *Agree or Disagree* | Do you agree or disagree with the following statement? People should be encouraged to take risks, even if there is a chance of failure. Use specific reasons and examples to support your answer. |
| *Preference* | Some people prefer to spend their money on experiences, such as travel or concerts, while others prefer to save for physical possessions, such as a car or a home. Which approach do you prefer, and why? Use specific reasons and examples to support your choice. |
| *If/Imaginary Situations* | If you could have any superpower, such as the ability to fly or become invisible, which one would you choose, and why? Use specific reasons and examples to explain your answer. |
| *Advan. and Disadvan.* | What are the advantages and disadvantages of living in a large city? Use specific reasons and examples to support your answer. |
| *Descriptive* | Describe a memorable trip you have taken and explain what made it special. Use specific details to support your response. |

Table 4: Examples of different question types and corresponding essay statements (prompts).

| Label | Train | Valid | Dev-Test | Total |
|---|---|---|---|---|
| **English** | | | | |
| *AI* | 925 | 299 | 712 | 1,936 |
| *Human* | 1,145 | 182 | 174 | 1,501 |
| **Total** | 2,070 | 481 | 886 | 3,437 |
| **Arabic** | | | | |
| *AI* | 1,467 | 391 | 369 | 2,227 |
| *Human* | 629 | 1,235 | 500 | 2,364 |
| **Total** | 2096 | 1,626 | 869 | 4,591 |

Table 5: Development phase: dataset and label distribution

to carry out this task. Each annotator was provided with clear anonymization guidelines and examples to ensure consistency and accuracy. Such anonymization steps ensure that the dataset meets ethical standards for research.

### 3.4.2 Data Generation

For the AI-generated essays, we followed two distinct methodologies:

- *Freehand Generation:* An instruct-tuned LLM, namely gpt-4o, independently generated essays using the *Freehand Generation Prompt* shown in Table 2. The prompt was de-

signed to ensure diverse outputs. We were inspired by the prompting techniques proposed by Chen et al. (2024).

- *Paraphrasing Human-Written Text:* Using the *Paraphrasing Prompt* shown in Table 3, human-authored essays were rephrased by an instruct-tuned LLM, namely claude-3.5 to generate stylistically varied yet semantically equivalent AI-written versions. The resulting text comprises a mix of human-written and AI-generated content, designed to challenge the effectiveness of detection methods.

| Category | English | Arabic | Total |
|---|---|---|---|
| AI (Free) | 400 | 100 | 500 |
| AI (Para) | 365 | 98 | 463 |
| Human | 365 | 95 | 460 |
| **Total** | 1,130 | 293 | 1,423 |

Table 6: Distribution of essays by *category* and *language* across the test set. Free - freehand generation, Para - paraphrasing-based generation.

The final GRACE dataset comprises a balanced distribution of human-written and AI-generated essays. Table 6 provides a detailed breakdown across languages and generation methods.

### 3.5 Baseline and Evaluation Setup

#### 3.5.1 Baseline

For all languages, we train an n-gram (unigram, $n = 1$) based baseline model. We transformed the texual content of the essays into a TF-IDF (Term Frequency-Inverse Document Frequency) representation with a maximum of 10k features. A Support Vector Machine (SVM) classifier is then trained on this feature representation to evaluate its performance.

#### 3.5.2 Evaluation Setup

The task was organized into two phases, corresponding to the previously described dataset development process:

- **Development phase**: We released the train and validation subsets, and participants submitted runs on the dev-test set through a competition on Codalab.[10]

- **Evaluation phase**: We released the official test subset – GRACE, and the participants were given four days to submit their final predictions through the same Codalab competition URL. Only the latest submission from each team was considered official and was used for the final team ranking.

#### 3.5.3 Evaluation Measure:

We measure the performance of the participating systems using accuracy, macro- precision, recall and F1 measure. However, official ranking was based on macro-F1.

## 4 Results and Overview of the Systems

In Table 7, we present the results of participants' systems for both Arabic and English including baseline. For Arabic, all systems outperformed the n-gram baseline, whereas, for English, three teams performed below the baseline. The task generated significant interest, with 56 teams registering to participate. However, the number of system submissions was nearly halved, and ultimately, only five teams submitted system description papers. In Table 8, we provide an overview of the participating systems for which a description paper was submitted. For Arabic top team, **IntegrityAI** (AL-Smadi, 2025), fine-tuned Electra model. For English top team, **CMI-AIGCX** (Kaijie et al., 2025), used LLMs (Llama 2 and 3) and also fine-tuned XLM-roberta model.

---

[10]https://codalab.lisn.upsaclay.fr/competitions/20118

Team **IntegrityAI** (AL-Smadi, 2025) fine-tuned ELECTRA-small for English and AraELECTRA-base for Arabic to balance high performance with computational efficiency. Stylometric features, including word count, sentence length, and vocabulary richness, were incorporated to enhance detection capabilities. The lightweight models achieved F1-scores of 0.985 for English and 0.984 for Arabic, demonstrating the effectiveness of combining transformer-based architectures with stylometric analysis. The system was further optimized for deployment on GPUs with moderate memory capacity, ensuring both efficiency and accessibility. Larger models, such as ELECTRA-large, were also tested, achieving an F1-score of 0.997 for English, demonstrating the potential for even greater accuracy with additional computational resources.

Team **CMI-AIGCX** (Kaijie et al., 2025) proposed a method leveraging the Llama-3.1-8B model as a proxy to capture the semantic feature of each token in the text. These token representations were subsequently used to train a model. Instead of fine-tuning an LLM, they leveraged multilingual knowledge and trained a model to enhance detection performance. Their approach demonstrated that using a proxy model with diverse multilingual knowledge can effectively detect machine-generated text across multiple languages, regardless of model size. For English, an F1 score of 0.999 was achieved, securing first place out of 25 teams. For Arabic, an F1 score of 0.965 was obtained, which ranked fourth among 21 teams.

Team **Tesla** (Indurthi and Varma, 2025) extracted a comprehensive set of features encompassing style, language complexity, bias, subjectivity, and emotion. These features were used to train four machine learning algorithms: Logistic Regression, Random Forest, Randomized Decision Trees (Extra Trees), and XGBoost, leveraging diverse approaches to optimize detection performance. Their methods ranked 6th on the leaderboard for the English subtask, achieving an F1-score of 0.986.

Team **EssayDetect** (Agrahari et al., 2025) proposed a fusion model by integrating pre-trained language model embeddings with stylometric and linguistic features to improve classification accuracy. The contributions were threefold: *(i)* LIME was utilized to identify and highlight highly discriminative features, *(ii)* focal loss was employed to address class imbalance, and *(iii)* layer-wise freezing was implemented during fine-tuning to preserve core linguistic representations in the lower layers while

## Table 7

**Arabic**

| Team | Acc | P | R | F1 | Rank |
|---|---|---|---|---|---|
| IntegrityAI | 0.986 | 0.990 | 0.979 | 0.984 | 1 |
| USTC-BUPT | 0.976 | 0.983 | 0.963 | 0.972 | 2 |
| starlight | 0.969 | 0.964 | 0.966 | 0.965 | 3 |
| CMI-AIGCX | 0.969 | 0.966 | 0.964 | 0.965 | 4 |
| apricity | 0.966 | 0.969 | 0.953 | 0.960 | 5 |
| RA | 0.962 | 0.956 | 0.959 | 0.957 | 6 |
| 1-800 | 0.959 | 0.961 | 0.945 | 0.952 | 7 |
| Lkminnow | 0.956 | 0.943 | 0.959 | 0.950 | 8 |
| alpaca0000001 | 0.949 | 0.937 | 0.948 | 0.942 | 9 |
| jojoc | 0.949 | 0.939 | 0.946 | 0.942 | 10 |
| small | 0.945 | 0.938 | 0.938 | 0.938 | 11 |
| jebish7 | 0.945 | 0.945 | 0.929 | 0.937 | 12 |
| EssayDetect | 0.942 | 0.949 | 0.919 | 0.932 | 13 |
| nits_teja_srikar | 0.922 | 0.943 | 0.882 | 0.904 | 14 |
| Mashixuan | 0.898 | 0.877 | 0.911 | 0.889 | 15 |
| Sinai | 0.829 | 0.821 | 0.866 | 0.822 | 16 |
| Vasudha | 0.816 | 0.796 | 0.831 | 0.804 | 17 |
| ShixuanMa | 0.758 | 0.783 | 0.818 | 0.754 | 18 |
| gaoyf | 0.608 | 0.720 | 0.707 | 0.607 | 19 |
| CNLP-NITS-PP | 0.590 | 0.557 | 0.563 | 0.557 | 20 |
| halcyonized | 0.495 | 0.488 | 0.487 | 0.475 | 21 |
| *Baseline* | 0.474 | 0.480 | 0.477 | 0.461 | - |

**English**

| Team | Acc | P | R | F1 | Rank |
|---|---|---|---|---|---|
| CMI-AIGCX | 0.999 | 0.999 | 0.999 | 0.999 | 1 |
| starlight | 0.997 | 0.998 | 0.996 | 0.997 | 2 |
| saehyunMa | 0.994 | 0.995 | 0.990 | 0.993 | 3 |
| Fsf | 0.994 | 0.995 | 0.990 | 0.993 | 4 |
| 1-800 | 0.991 | 0.987 | 0.993 | 0.990 | 5 |
| Tesla | 0.988 | 0.983 | 0.989 | 0.986 | 6 |
| apricity | 0.988 | 0.983 | 0.989 | 0.986 | 7 |
| small | 0.984 | 0.981 | 0.983 | 0.982 | 8 |
| jojoc | 0.982 | 0.975 | 0.985 | 0.980 | 9 |
| EssayDetect | 0.978 | 0.968 | 0.984 | 0.975 | 10 |
| ShixuanMa | 0.976 | 0.968 | 0.979 | 0.973 | 11 |
| RA | 0.973 | 0.975 | 0.964 | 0.969 | 12 |
| alpaca0000001 | 0.956 | 0.940 | 0.967 | 0.951 | 13 |
| Lkminnow | 0.932 | 0.913 | 0.943 | 0.925 | 14 |
| IntegrityAI | 0.880 | 0.864 | 0.911 | 0.873 | 15 |
| USTC-BUPT | 0.878 | 0.922 | 0.812 | 0.842 | 16 |
| jebish7 | 0.847 | 0.908 | 0.763 | 0.794 | 17 |
| CNLP-NITS-PP | 0.777 | 0.784 | 0.825 | 0.771 | 18 |
| Mashixuan | 0.742 | 0.778 | 0.809 | 0.739 | 19 |
| nits_teja_srikar | 0.773 | 0.875 | 0.649 | 0.658 | 20 |
| Vasudha | 0.517 | 0.700 | 0.643 | 0.509 | 21 |
| Mahavir_IIITA | 0.512 | 0.683 | 0.634 | 0.504 | 22 |
| *Baseline* | 0.495 | 0.494 | 0.494 | 0.478 | - |
| halcyonized | 0.493 | 0.494 | 0.493 | 0.477 | 23 |
| gaoyf | 0.391 | 0.523 | 0.514 | 0.374 | 24 |
| Sinai | 0.354 | 0.602 | 0.519 | 0.298 | 25 |

Table 7: The official results for Arabic and English are ranked based on the official metric: macro-F1. Teams that submitted a system description paper are indicated in bold.

## Table 8

| Team | Lang.<br>Arabic | English | LLama2 | LLama3 | BERT | RoBERTa | XLM-r | ALBERT | DistilBERT | DeBERTa | Electra | AraBERT | Prep. | Info. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| IntegrityAI | 1 | 15 | | | | | | | | | | ☑ | ☑ | ☑ |
| CMI-AIGCX | 4 | 1 | ☑ | ☑ | | | ☑ | | | | | | ☑ | |
| Tesla | | 6 | | | | | | | | | | | | |
| EssayDetect | 13 | 10 | | | ☑ | ☑ | ☑ | ☑ | ☑ | | | | | |
| RA | 6 | 12 | | | | ☑ | | | | ☑ | | ☑ | | |

Table 8: Overview of the approaches. The numbers in the language box refer to the position of the team in the official ranking. Prep.: Preprocessing. Info.: Info. Extraction.

enabling the higher layers to capture task-specific stylistic differences in essays.

Team **RA** (Gharib and Elgendy, 2025) fine-tuned several models for English, including RoBERTa, XLM-RoBERTa, mBERT, and DeBERTa. Similar performance was observed across all models on the validation set, except for mBERT, which exhibited slightly lower performance. For Arabic, AraBERT, ArBERT, and MarBERT were fine-tuned on the full dataset. AraBERT consistently demonstrated superior performance in terms of F1-score across both languages. The models consistently exceeded both the mean and median scores across tasks, achieving an F1-score of 0.969 in classifying AI-generated essays in English and 0.957 in Arabic.

## 5 Conclusion and Future Work

We presented an overview of the shared task on the *Academic Essay Challenge*. The task attracted significant attention, with a total of 56 teams registering to participate in the development and evaluation phases. Of these, 21 teams submitted official results on the test set for Arabic, and 25 teams did so for English. Finally, seven teams submitted task description papers. Most systems fine-tuned transformer-based language models; however, several teams also incorporated additional features, such as style, language complexity, bias, subjectivity, and emotion. For both languages, the top-performing teams achieved F1 scores above 0.98.

## Limitations

A major limitation of the dataset is its small size, particularly for Arabic, which restricts the development of more robust models. The challenging nature of academic essay collection is reflected in the limited dataset size. Future studies could focus on curating larger datasets to enable the creation of more challenging tasks and the development of more robust models.

## Ethical Considerations

The datasets used in the shared task may reflect subjective biases or perspectives of the essay authors, even though they followed the provided instructions. Importantly, the datasets do not include any personal information, and no such information was collected during the data curation process. Therefore, we do not anticipate any ethical concerns related to privacy. Furthermore, the dataset was shared only with participants who signed an agreement, ensuring responsible use of the dataset.

## Acknowledgments

## References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.

Shifali Agrahari, Subhashi Jayant, Saurabh Kumar, and Ranbir Sanasam. 2025. Team EssayDetect at GenAI Detection Task 2: Guardians of Academic Integrity: Multilingual Detection of AI-Generated Essays. In *Proceedings of the 1st Workshop on GenAI Content Detection (GenAIDetect)*, Abu Dhabi, UAE. International Conference on Computational Linguistics.

Abdelhamid M. Ahmed, Xiao Zhang, Lameya M. Rezk, and Wajdi Zaghouani. 2023. Building an annotated l1 arabic/l2 english bilingual writer corpus: The qatari corpus of argumentative writing (qcaw). *Corpus-based Studies across Humanities*, 1(1):183–215.

Mohammad AL-Smadi. 2025. IntegrityAI at GenAI Detection Task 2: Detecting Machine-Generated Academic Essays in English and Arabic Using ELECTRA and Stylometry. In *Proceedings of the 1st Workshop on GenAI Content Detection (GenAIDetect)*, Abu Dhabi, UAE. International Conference on Computational Linguistics.

AYG Alfaifi and ES Atwell. 2013. Arabic learner corpus v1: A new resource for arabic language research. In *Second Workshop on Arabic Corpus Linguistics*.

Amrita Bhattacharjee, Tharindu Kumarage, Raha Moraffah, and Huan Liu. 2023. ConDA: Contrastive domain adaptation for AI-generated text detection. In *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 598–610, Nusa Dua, Bali. Association for Computational Linguistics.

Jiuhai Chen, Rifaa Qadri, Yuxin Wen, Neel Jain, John Kirchenbauer, Tianyi Zhou, and Tom Goldstein. 2024. Genqa: Generating millions of instructions from a handful of prompts. *Preprint*, arXiv:2406.10323.

Kohinoor Monish Darda, Marion Carre, and Emily S. Cross. 2023. Value attributed to text-based archives generated by artificial intelligence. *Royal Society Open Science*, 10(2):220915.

Liam Dugan, Alyssa Hwang, Filip Trhlík, Andrew Zhu, Josh Magnus Ludan, Hainiu Xu, Daphne Ippolito, and Chris Callison-Burch. 2024. RAID: A shared benchmark for robust evaluation of machine-generated text detectors. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12463–12492, Bangkok, Thailand. Association for Computational Linguistics.

Matthias Gallé, Jos Rozen, Germán Kruszewski, and Hady Elsahar. 2021. Unsupervised and distributional detection of machine-generated text. *arXiv preprint arXiv:2111.02878*.

Rana Gharib and Ahmed Elgendy. 2025. RA at GenAI Detection Task 2: Fine-tuned Language Models For Detection of Academic Authenticity, Results and Thoughts. In *Proceedings of the 1st Workshop on GenAI Content Detection (GenAIDetect)*, Abu Dhabi, UAE. International Conference on Computational Linguistics.

Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How close is chatgpt to human experts? comparison corpus, evaluation, and detection. *arXiv preprint arXiv:2301.07597*.

Hans WA Hanley and Zakir Durumeric. 2024. Machine-made media: Monitoring the mobilization of machine-generated articles on misinformation and mainstream news websites. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 18, pages 542–556.

Xuanli He, Qiongkai Xu, Lingjuan Lyu, Fangzhao Wu, and Chenguang Wang. 2022. Protecting intellectual property of language generation apis with lexical

watermark. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10):10758–10766.

Vijayasaradhi Indurthi and Vasudeva Varma. 2025. Tesla at GenAI Content Detection Task 1: LLM Agents in Multilingual Machine-Generated Text Detection. In *Proceedings of the 1st Workshop on GenAI Content Detection (GenAIDetect)*, Abu Dhabi, UAE. International Conference on Computational Linguistics.

Jiao Kaijie, Yao Xingyu, Ma Shixuan, Fang Sifan, Guo Zikang, Xu Benfeng, Zhang Licheng, Wang Quan, Zhang Yongdong, and Mao Zhendong. 2025. CMI-AIGCX at GenAI Detection Task 2: Leveraging Multilingual Proxy LLMs for Machine-Generated Text Detection in Academic Essays. In *Proceedings of the 1st Workshop on GenAI Content Detection (GenAIDetect)*, Abu Dhabi, UAE. International Conference on Computational Linguistics.

John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. A watermark for large language models. In *International Conference on Machine Learning*, pages 17061–17084. PMLR.

Weixin Liang, Mert Yuksekgonul, Yining Mao, Eric Wu, and James Zou. 2023. GPT detectors are biased against non-native english writers. *Patterns*, 4(7):100779.

George Mikros, Athanasios Koursaris, Dimitrios Bilianos, and George Markopoulos. 2023. AI-writing detection using an ensemble of transformers and stylometric features. In *Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2023) co-located with the Conference of the Spanish Society for Natural Language Processing (SEPLN 2023)*, volume 3496 of *CEUR Workshop Proceedings*, pages 1–14, Jaén, Spain.

OpenAI. 2024. Gpt-4o mini: Advancing cost-efficient intelligence. *OpenAI Blog*.

Mike Perkins, Jasper Roe, Binh H. Vu, Darius Postma, Don Hickerson, James McGaughran, and Huy Q. Khuat. 2024. Simple techniques to bypass GenAI text detectors: Implications for inclusive education. *International Journal of Educational Technology in Higher Education*, 21(1):53.

Hooman H. Rashidi, Brandon D. Fennell, Samer Albahra, Bo Hu, and Tom Gorbett. 2023. The chatgpt conundrum: Human-generated scientific manuscripts misidentified as ai creations by ai text detection tool. *Journal of Pathology Informatics*, 14:100342.

Aditya Shah, Prateek Ranka, Urmi Dedhia, Shruti Prasad, Siddhi Muni, and Kiran Bhowmick. 2023. Detecting and unmasking AI-generated texts through explainable artificial intelligence using stylistic features. *International Journal of Advanced Computer Science and Applications*, 14(10):1–10.

Sebastian Szyller, Buse Gul Atli, Samuel Marchal, and N. Asokan. 2021. Dawn: Dynamic adversarial watermarking of neural networks. In *Proceedings of the 29th ACM International Conference on Multimedia*, MM '21, page 4417–4425, New York, NY, USA. Association for Computing Machinery.

Ruixiang Tang, Yu-Neng Chuang, and Xia Hu. 2024. The science of detecting llm-generated text. *Communications of the ACM*, 67(4):50–59.

Gemini Team. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.

Vismay Vora, Jenil Savla, Deevya Mehta, and Aruna Gawade. 2023. A multimodal approach for detecting AI generated content using BERT and CNN. *International Journal on Recent and Innovation Trends in Computing and Communication*, 11(9):691–701.

Yuxia Wang, Jonibek Mansurov, Petar Ivanov, Jinyan Su, Artem Shelmanov, Akim Tsvigun, Osama Mohammed Afzal, Tarek Mahmoud, Giovanni Puccetti, Thomas Arnold, Thomas Arnold, Chenxi Whitehouse, Alham Fikri Aji, Nizar Habash, Iryna Gurevych, and Preslav Nakov. 2024a. SemEval-2024 Task 8: Multidomain, multimodel and multilingual machine-generated text detection. In *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*, pages 2057–2079, Mexico City, Mexico. Association for Computational Linguistics.

Yuxia Wang, Jonibek Mansurov, Petar Ivanov, Jinyan Su, Artem Shelmanov, Akim Tsvigun, Chenxi Whitehouse, Osama Mohammed Afzal, Tarek Mahmoud, Toru Sasaki, Thomas Arnold, Alham Aji, Nizar Habash, Iryna Gurevych, and Preslav Nakov. 2024b. M4: Multi-generator, multi-domain, and multilingual black-box machine-generated text detection. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1369–1407, St. Julian's, Malta. Association for Computational Linguistics.

Debora Weber-Wulff, Alla Anohina-Naumeca, Sonja Bjelobaba, Tomáš Foltýnek, Jean Guerrero-Dib, Olumide Popoola, Petr Šigut, and Lorna Waddington. 2023. Testing of detection tools for AI-generated text. *International Journal for Educational Integrity*, 19(1):1–39.

Hongyu Wu and Tom Flanagan. 2023. The limits of AI content detectors. *Journal of Student Research*, 12(3):1–7.

Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Derek F Wong, and Lidia S Chao. 2023. A survey on llm-gernerated text detection: Necessity, methods, and future directions. *arXiv preprint arXiv:2310.14724*.

Peipeng Yu, Jiahan Chen, Xuan Feng, and Zhihua Xia. 2023. Cheat: A large-scale dataset for detecting chatgpt-written abstracts. *arXiv preprint arXiv:2304.12008*.

Wajdi Zaghouani, Abdelhamid Ahmed, Xiao Zhang, and Lameya Rezk. 2024. QCAW 1.0: Building a qatari corpus of student argumentative writing. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 13382–13394, Torino, Italia. ELRA and ICCL.

Wataru Zaitsu and Mingzhe Jin. 2023. Distinguishing chatgpt(-3.5, -4)-generated and human-written papers through japanese stylometric analysis. *PLOS ONE*, 18(8):1–12.