

# The State of the Art in Detecting LLM-Generated Text in Academic Assignments

Author: Manus AI

Date: November 22, 2025

## Introduction

The rapid proliferation of Large Language Models (LLMs) like ChatGPT has introduced significant challenges to academic integrity. As these models become increasingly sophisticated, their ability to generate human-like text makes it difficult for educators to distinguish between original student work and AI-generated content. This report provides a comprehensive overview of the current state-of-the-art in LLM-generated text detection, with a particular focus on the techniques and methods being used in academic settings. We will explore the primary detection paradigms, analyse their effectiveness, and discuss the significant challenges and limitations that currently exist.

## Primary Detection Paradigms

The detection of LLM-generated text is fundamentally a binary classification task: determining whether a given piece of text was written by a human or generated by a machine. Research in this field has coalesced around four primary paradigms, each with its own set of methodologies and challenges <sup>[1]</sup>.

Detection Paradigm	Description	Key Methods
Watermarking	Proactively embedding a hidden, statistical signal into the generated text that can be later detected.	Modifying the token selection process during generation (e.g., Tournament Sampling).
Statistical Analysis	Analysing the statistical properties of the text to identify patterns characteristic of machine generation.	Measuring features like perplexity (randomness) and burstiness (variation in perplexity).
Neural Network Classifiers	Training deep learning models to distinguish between human and AI-generated text based on complex linguistic patterns.	Fine-tuning transformer models like BERT and RoBERTa on large datasets of human and AI text.
Human-Assisted Methods	Relying on human experts to identify subtle cues and inconsistencies that automated systems might miss.	Manual review, stylistic analysis, and checking for factual inaccuracies or nonsensical statements.

These approaches are not mutually exclusive and are often used in combination to improve detection accuracy. However, as we will explore, each method faces significant limitations, particularly as LLMs become more advanced and users develop sophisticated evasion techniques.

- [1] Wu, J., et al. (2025). A Survey on LLM-Generated Text Detection: Necessity, Methods, and Future Directions. *\*Computational Linguistics\**, 51(1), 275–338.

### Watermarking Techniques

Watermarking represents a proactive approach to LLM detection. Instead of retrospectively analysing a finished piece of text, watermarking involves embedding an imperceptible, cryptographic signal into the text as it is being generated. This signal can then be detected later to confirm the text's origin.

One of the most promising and production-ready watermarking schemes is SynthID-Text, developed by Google DeepMind [2]. This method works by modifying the LLM's token sampling procedure. During text generation, it subtly biases the selection of tokens towards those that score highly under a set of secret, random watermarking functions. This creates a statistical pattern that is invisible to the human eye but can be reliably detected by an algorithm that knows the secret key. The detection process involves calculating a score based on how well the text aligns with these hidden functions. The strength of the watermark signal depends on two key factors: the length of the text and the entropy (randomness) of the LLM's output. Longer texts provide more evidence, and higher entropy allows for more flexibility in embedding the watermark without degrading text quality.

*"SynthID-Text does not affect LLM training and modifies only the sampling procedure; watermark detection is computationally efficient, without using the underlying LLM." [2]*

While powerful, watermarking is not a silver bullet. Its effectiveness is reduced when the LLM's output is very predictable (low entropy), as there is less room to embed the signal. Furthermore, watermarking requires the cooperation of the LLM provider to implement the technique at the source. This means that it is not a solution for detecting text from open-source or proprietary models that do not have watermarking enabled.

- [2] Dathathri, S., et al. (2024). Scalable watermarking for identifying large language model outputs. *\*Nature\**, 634, 818-823.

### Statistical Analysis

Statistical methods analyse the mathematical properties of a text to find signatures of machine generation. These methods are computationally efficient and form the basis of many commercially available AI detectors. The two most prominent statistical features used for this purpose are perplexity and burstiness [3].

Perplexity measures the randomness or unpredictability of a text. It evaluates how well a language model can predict the next word in a sequence. Human writing tends to be more unpredictable and varied, resulting in higher perplexity scores. In contrast, LLMs often choose the most probable and predictable words, leading to text with lower perplexity. A high perplexity score suggests that the text is less likely to have been generated by an AI.

Burstiness measures the variation in perplexity across a document. Humans write in bursts of creativity and inspiration, leading to significant fluctuations in sentence length, structure, and word choice. This results in a high burstiness score. LLMs, on the other hand, tend to be more formulaic

and consistent in their output, producing text with a uniform level of perplexity and therefore low burstiness. A low burstiness score is a strong indicator of AI-generated content.

*“While a person could easily write an AI-like sentence by accident, people tend to vary their sentence construction and diction throughout a document. On the other hand, models formulaically use the same rule to choose the next word in the sentence, leading to low burstiness.” [3]*

While these statistical measures can be effective, they are not foolproof. As LLMs become more sophisticated, they are better able to mimic the statistical properties of human writing. Furthermore, simple paraphrasing or manual editing of AI-generated text can significantly alter its perplexity and burstiness, making it much harder to detect.

[3] Tian, E. (2023, March 1). What is perplexity & burstiness for AI detection? \*GPTZero\*.

### Neural Network Classifiers

Neural network classifiers represent a more powerful, but also more computationally expensive, approach to AI detection. These methods utilize deep learning models, particularly transformer-based architectures like BERT (Bidirectional Encoder Representations from Transformers), to learn the complex linguistic patterns that differentiate human and AI-generated text [3]. By training on vast datasets containing examples of both, these models can identify subtle nuances in syntax, semantics, and style that are often missed by statistical methods.

However, the effectiveness of these classifiers is highly dependent on their training data. A landmark study from the University of Pennsylvania, which developed the Robust AI Detector (RAID) benchmark, found that most neural classifiers have critical vulnerabilities [4]. The study revealed two major weaknesses:

1. Cross-LLM Detection Failure: A model trained to detect output from one LLM (e.g., ChatGPT) is often “mostly useless” at detecting text from another (e.g., Llama).
2. Domain Specificity: A classifier trained on one genre of text, such as news articles, does not perform well when applied to another, like academic essays or creative writing.

This creates a constant “arms race” where detectors must be continually retrained to keep up with the latest LLMs. Furthermore, a highly influential 2023 study found that the accuracy of these tools dropped significantly when moving from detecting GPT-3.5 to the more advanced GPT-4 [6]. The same study also highlighted the persistent problem of false positives, where human-written text is incorrectly flagged as AI-generated, a critical issue in academic settings where such an accusation has serious consequences.

The unreliability of these systems is so pronounced that even OpenAI, the creator of ChatGPT, shut down its own AI detection tool in July 2023 due to its “low rate of accuracy” [5]. This underscores the current consensus that neural classifiers, while powerful in specific contexts, are not yet robust enough for widespread, high-stakes deployment in universities.

[4] Callison-Burch, C., & Dugan, L. (2024, August 12). Detecting Machine-Generated Text: An Arms Race With the Advancements of Large Language Models. \*Penn Engineering Blog\*.

[5] MIT Sloan Teaching & Learning Technologies. (n.d.). AI Detectors Don’t Work. Here’s What to Do Instead.

[6] Elkhataf, A. M., et al. (2023). Evaluating the efficacy of AI content detection tools in differentiating between human and AI-generated text. *\*International Journal for Educational Integrity\**, 19, 17.

## Challenges and Limitations

The arms race between LLM development and detection has created a landscape fraught with challenges. The most significant limitation of current detection tools is their declining accuracy in the face of more advanced models and simple evasion techniques.

### The Evasion Problem

A critical vulnerability of nearly all detection methods is their susceptibility to adversarial attacks. The Penn Engineering study found that simple manipulations can render most detectors useless. These include:

- Paraphrasing: Manually or automatically rephrasing AI-generated text.
- Inserting extra spaces or typos.
- Swapping letters for look-alike symbols (e.g., using a zero for the letter 'o').
- Using synonyms or alternative spellings.

These simple edits are often enough to break the statistical patterns and linguistic signatures that detectors rely on, allowing AI-generated text to pass as human-written.

### The False Positive Crisis

Perhaps the most pressing issue for universities is the high rate of false positives. Multiple studies have confirmed that detection tools frequently misclassify human-written text as AI-generated <sup>[5]</sup> <sup>[6]</sup>. This is particularly true for text written by non-native English speakers, whose writing style can sometimes mimic the simpler sentence structures that detectors associate with AI. The potential for falsely accusing a student of academic misconduct based on a faulty algorithm has led many universities to abandon or discourage the use of these tools altogether.

*"AI detection software is far from foolproof – in fact, it has high error rates and can lead instructors to falsely accuse students of misconduct." [5]*

## Conclusion

The state of the art in LLM-generated text detection is a dynamic and rapidly evolving field. While techniques like watermarking, statistical analysis, and neural network classifiers show promise, they are currently outpaced by the advancement of LLMs and the ease with which their outputs can be disguised. The high error rates and susceptibility to evasion of current detection tools make them unreliable for high-stakes academic assessments. The consensus among researchers and many academic institutions is that a purely technological solution to this challenge is not yet viable.

Instead, the focus is shifting towards pedagogical approaches. This includes designing AI-resistant assignments that require critical thinking, personal reflection, and in-class participation, as well as fostering a culture of academic integrity through clear policies and open dialogue with students. While the technological arms race will undoubtedly continue, the most effective strategy for universities in the near term lies in adapting their teaching and assessment methods to the new reality of AI-assisted writing.