

---

# What Are We Detecting, Really? LLM-Generated Text Detection Remains an Unsolved Problem

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This position paper argues that, in most practical cases, it is not possible to accurately  
2 detect LLM-generated text. We consider that “LLM-generated text” refers  
3 to the content produced by LLMs through normal prompts. As implied by the  
4 names “LLM-generated text” and “human-written text”, the difference lies in how  
5 they are produced, but in practice, we can only evaluate them based on the final  
6 output—the text itself—where there is often significant overlap between human-  
7 and machine-generated content.. The numerical results of LLM-generated text de-  
8 tection are often misunderstood and their significance is diminishing. The detectors  
9 can serve a purpose under specific conditions, whose results should only be used  
10 as a reference with greater caution rather than the decisive indicator.

## 11 1 Introduction

12 The rapid development of large language models (LLMs) has led to a rise of LLM-generated text,  
13 which has been observed across various fields, such as academia [Liang et al., 2024, Geng and  
14 Trotta, 2024], Wikipedia [Brooks et al., 2024, Huang et al., 2025], and numerous online texts [Sun  
15 et al., 2024, Liang et al., 2025]. The detection of LLM-generated text has attracted the attention of  
16 researchers, and many detectors have been proposed and studied [Yang et al., 2023, Wu et al., 2025].  
17 Before starting our discussion, we want to clarify the following definition:

18 *What exactly is “LLM-generated text”?*

19 In fact, the term “LLM-generated text” is fairly new. Researchers also used expressions like “machine-  
20 generated text” or “AI-generated”, as seen in Table 1. For simplicity, we use “LLM-generated text”  
21 to refer to the subject of study in this paper, as it is more precise than the other expressions. Apart  
22 from slight differences in terminology, the definition of “LLM-generated text” in most papers is quite  
23 broad, meaning the text can be produced in many ways using LLMs, like paraphrasing, translation, or  
24 generating long text from simple prompts. Therefore, we can consider that **“LLM-generated text”**  
25 **refers to the content produced by LLMs through normal prompts.**

26 The reliability of the detectors has also been widely discussed [Sadasivan et al., 2023, Chakraborty  
27 et al., 2024]. The indistinguishability between LLM-generated and human-written text is one of big  
28 challenges for LLMs [Kaddour et al., 2023]. Similar to the central question in Chakraborty et al.  
29 [2024]’s work, we explore the following key one:

30 *Is it possible to detect the LLM-generated text in practice?*

31 Literally speaking, LLM-generated text detectors need to cover all these different scenarios, but few  
32 detectors have tried to distinguish them [Cheng et al., 2025]. If we take into account the diversity of  
33 LLMs and prompts, as well as human-in-the-loop, the situation becomes even more intricate. Hence,  
34 we argue that, **in most practical cases, it is not possible to accurately detect LLM-generated text.**

Table 1: Definition of LLM-generated text in different papers

Paper	Definition
Crothers et al. [2023]	<i>“Machine-generated text” is natural language text that is produced, modified, or extended by a machine.</i>
Kumarage et al. [2024]	<i>In this survey, we define <u>AI-generated text</u> as output produced by a natural language generation pipeline employing a neural probabilistic language model.</i>
Wu et al. [2025]	<i>LLM-generated Text is defined as cohesive, grammatically sound, and pertinent content generated by LLMs.</i>

35 While limitations of these detection methods have caused concern among researchers [Sadasivan  
 36 et al., 2023, Liang et al., 2023, Doughman et al., 2024, Nicks et al., 2023, Saha and Feizi, 2025], they  
 37 could be applied in diverse contexts. The emergence of the GPTZero platform is a good example,  
 38 although it is currently unclear how frequently people use it. So the next question arises naturally:

39 *Should we use these detectors?*

40 We think that **the detectors can serve a purpose under specific conditions, whose results should**  
 41 **only be used as a reference with greater caution rather than the decisive indicator.**

42 We will discuss them in detail in the following sections.

## 43 2 Detectors

44 Probably most people became aware of LLMs after the release of ChatGPT, but the research on  
 45 detecting text generated by language models had started before that. For example, Gehrmann et al.  
 46 [2019] proposed the GLTR tool to detect whether text was generated by models, with experiments  
 47 involving GPT-2 [Radford et al., 2019] and BERT [Devlin et al., 2019]. Zellers et al. [2019] developed  
 48 the Grover model to detect AI-generated fake news. Even GPT-3 [Brown et al., 2020] continued to  
 49 face skepticism regarding its text-generation capabilities [Bender et al., 2021], making the detectors’  
 50 performance unsurprising.

51 Another pioneering work by [Ippolito et al., 2019] demonstrated that humans have already encountered  
 52 some difficulties in identifying texts generated by GPT-2. Later, Clark et al. [2021] found that  
 53 untrained people at the time were not very good at recognizing text produced by GPT-3, and Wahle  
 54 et al. [2022] noticed the similar situation for machine-paraphrased plagiarism.

55 In the past two or three years, the rapid development and spread of LLMs has drawn significant  
 56 attention from researchers to the detection of LLM-generated text, and diverse methods have been  
 57 proposed [Wu et al., 2025]: DetectGPT [Mitchell et al., 2023], Fast-DetectGPT [Bao et al., 2023],  
 58 DetectLLM [Su et al., 2023], LLMDet [Wu et al., 2023], DeID-GPT [Liu et al., 2023] and some  
 59 others [Dugan et al., 2023] in 2023; Binoculars [Hans et al., 2024], TOCSIN [Ma and Wang, 2024],  
 60 Dpic [Yu et al., 2024b], Text Fluoroscopy [Yu et al., 2024a] in 2024. The examples listed above are  
 61 illustrative, and the actual number of detectors is much larger.

62 In the meantime, specialized detectors have been developed, for instance, targeting journalistic  
 63 news articles [Bhattacharjee et al., 2023] and tweets [Kumarage et al., 2023, Gambini et al., 2022].  
 64 Additionally, the detection of LLM-generated text is not limited to English [Wang et al., 2025]. we  
 65 have also seen detectors for other languages, such as French [Antoun et al., 2023a], Japanese [Zaitsu  
 66 and Jin, 2023], Chinese [Wang et al., 2024a].

67 While these techniques of detection performed well earlier on certain datasets, the ongoing progress  
 68 of LLMs also makes detection harder [Wu et al., 2025]. A wide range of methods are utilized by  
 69 these detectors, but the absence of universal benchmarks and different application scenarios limit  
 70 meaningful comparison. We will address this issue in greater detail in Section 5.1.

71 These detection methods can be classified into many categories according to different criteria. For  
 72 instance, Abdali et al. [2024] classifies them as supervised methods, zero-shot methods, retrieval-  
 73 based methods, watermarking methods, discriminating features. Wu et al. [2025] mainly examines

74 them through the lens of watermarking techniques, statistics-based detectors, neural-based detectors,  
75 and human-assisted methods. It is difficult to provide a comprehensive summary of LLM-generated  
76 text detectors, but to our knowledge, no detector has been conclusively established as the best,  
77 particularly in practical deployment contexts.

78 There are other ways to categorize the detectors. For example, most studies only think about binary  
79 classification, and detectors with multi-category cases have rarely been explored, which will be  
80 further examined in Section 5.4.

### 81 3 Related Work

82 **Benchmark** One of the major challenges in establishing benchmarks for detecting LLM-generated  
83 text is that LLMs are continuously evolving, and their characteristics do not remain the same.  
84 For example, [Liyanage et al., 2022] created their benchmark with GPT-2, which should be quite  
85 differently from the current advanced LLMs. More LLMs were employed in subsequent benchmark  
86 construction [Wang et al., 2024b, He et al., 2024, Cornelius et al., 2024], but the number of prompts  
87 and scenarios used was limited. Some recent benchmarks [Tao et al., 2024, Wu et al., 2024] have  
88 incorporated a broader range of scenarios, and their impact and effectiveness remain to be seen.  
89 Similar challenge for the dataset [Gritsai et al., 2024].

90 **Watermarking and Attack** As mentioned earlier, watermarking LLMs is considered a category of  
91 detection methods, and it has shown good effectiveness in simulation [Kirchenbauer et al., 2023],  
92 which may also be an ethical necessity [Grinbaum and Adomaitis, 2022]. Researchers have proposed  
93 diverse watermarking techniques [Jovanović et al., 2024, Dathathri et al., 2024, Christ et al., 2024],  
94 but given the variety of attack techniques, this game is far from over [Krishna et al., 2023, Cai and  
95 Cui, 2023, Sadasivan et al., 2023, Lau and Zubiaga, 2024, Tufts et al., 2024, Macko et al., 2024].

96 **Evaluation and Test** These different methods have also been tested by various researchers [Weber-  
97 Wulff et al., 2023, Walters, 2023], in the beginning, some issues have been pointed out, such as  
98 biases [Liang et al., 2023] and false positives [Rashidi et al., 2023, Giray, 2024]. Although different  
99 detectors have been proposed over time, many challenges still exist [Kar et al., 2024, Tufts et al.,  
100 2024]. Given the diversity of usage scenarios and the variety of LLMs, such results are to be expected.  
101 We will provide a detailed analysis later.

### 102 4 Case Study

103 For a simple illustration of LLM-generated text and detection, we take the introduction from the  
104 website of NeurIPS<sup>1</sup> for simulations. This paragraph is processed with various prompts across  
105 different LLMs (DeepSeek-V3 [Liu et al., 2024], DeepSeek-R1 [Guo et al., 2025], GPT-3.5, GPT-4o-  
106 mini, GPT4o [Hurst et al., 2024]).

107 We present the original text and two generated texts with different prompts, and more examples are  
108 provided in the Appendix A.

109 **Original Text** *The conference was founded in 1987 and is now a multi-track interdisciplinary*  
110 *annual meeting that includes invited talks, demonstrations, symposia, and oral and poster presenta-*  
111 *tions of refereed papers. Along with the conference is a professional exposition focusing on machine*  
112 *learning in practice, a series of tutorials, and topical workshops that provide a less formal setting for*  
113 *the exchange of ideas.*

114 **DeepSeek-V3** (“Polish the following passage: ...”) *Founded in 1987, the conference has grown*  
115 *into a premier multi-track interdisciplinary annual event. It features invited talks, demonstrations,*  
116 *symposia, and peer-reviewed oral and poster presentations. Alongside the conference, attendees can*  
117 *explore a professional exposition highlighting practical applications of machine learning, a series of*  
118 *tutorials, and topical workshops that foster informal idea exchange.*

---

<sup>1</sup><https://neurips.cc/>

119 **DeepSeek-R1** (“Polish the following passage: ...”) *Established in 1987, the conference has*  
 120 *grown into a premier interdisciplinary event held annually, featuring multiple tracks of programming.*  
 121 *Attendees engage with invited talks, demonstrations, symposia, and peer-reviewed research presented*  
 122 *through both oral and poster sessions. Complementing the core program are a professional exposition*  
 123 *showcasing real-world applications of machine learning, hands-on tutorials, and topical workshops*  
 124 *designed to foster informal dialogue and collaborative idea exchange.*

125 **GPT-4o** (“Polish the following passage: ...”) *Established in 1987, the conference has evolved into*  
 126 *a comprehensive, multi-track interdisciplinary annual event. It features invited talks, demonstrations,*  
 127 *symposia, and both oral and poster presentations of peer-reviewed papers. In addition to the main*  
 128 *conference, there is a professional exposition dedicated to practical applications of machine learning,*  
 129 *as well as a series of tutorials and topical workshops. These workshops offer a more informal setting*  
 130 *for the exchange of ideas.*

131 **GPT-4o** (“Rewrite the following passage without using complex expressions: ...”) *The confer-*  
 132 *ence started in 1987 and now happens every year with multiple sessions. It features invited talks,*  
 133 *demonstrations, discussions, and presentations of selected papers. There is also a professional expo*  
 134 *about using machine learning, a series of tutorials, and workshops that offer a more relaxed space*  
 135 *for sharing ideas.*

136 Although they modify the original text (for example, the underlined words above), the added words  
 137 and expressions are not the same. Table 2 presents the results of detecting these texts using Fast-  
 138 DetectGPT [Bao et al., 2023]. Even though these texts are all generated by LLMs, their detection  
 139 outcomes vary widely. As we can easily find, in this case, the detector struggles to clearly identify text  
 140 generated by DeepSeek-V3 and DeepSeek-R1, as the probability of texts by them being identified as  
 141 machine-generated is even lower than that of the original text from NeurIPS website. The results of  
 142 GPT-3.5, GPT-4o-mini, and GPT-4o show that prompts can easily affect the outputs and the detection  
 143 results.

Table 2: Detection results using Fast-DetectGPT. The last two columns correspond to the predictions of the machine-generated results when the Sampling/scoring model is gpt-neo-2.7b and falcon-7b, respectively.

Prompts	Model	p1	p2
(original text)	-	44%	23%
Polish the following passage:	DeepSeek-V3	34%	12%
	DeepSeek-R1	<b>21%</b>	<b>10%</b>
Polish the following passage:	GPT-3.5	50%	22%
	GPT-4o-mini	35%	20%
	GPT-4o	<b>84%</b>	<b>75%</b>
Rewrite the following passage without using complex expressions:	GPT-3.5	30%	16%
	GPT-4o-mini	51%	55%
	GPT-4o	38%	15%

144 These are merely a few basic examples of the issues and limitations faced by LLM-generated text  
 145 detectors. A more detailed discussion will follow in the next section.

## 146 5 Issues and Limitations

147 As we briefly introduced before, the detection of text generated by LLMs has emerged as a widely  
 148 discussed and actively pursued task in natural language processing. Such detection tools are often  
 149 promoted for their potential utility in identifying instances of plagiarism [Pudasaini et al., 2024],  
 150 academic dishonesty (e.g., cheating during examinations) [Wang and Li, 2025], the automatic  
 151 generation of unethical peer reviews [Kumar et al., 2025], and other forms of content manipulation.

152 Meanwhile, many issues and challenges have been discussed [Tang et al., 2024, Wu et al., 2025,  
 153 Fraser et al., 2024, Abdali et al., 2024]. Despite their apparent usefulness, there are some fundamental

154 limitations associated with these tools that raise serious ethical and methodological concerns. We  
155 will address these issues and limitations from various perspectives in this section.

## 156 5.1 Lack of Precise Definition and Gold-Standard Benchmark

157 Unlike most question-answering or classification tasks, “human-written text” itself lacks a clear  
158 and well-defined boundary compared to “LLM-generated text”. **As implied by their names, the  
159 difference lies in how they are produced, but in practice, we can only assess them based on their  
160 final output i.e., the text, where in which a lot of overlap between them.**

161 Researchers often say that the text generated by LLMs is different from that written by hu-  
162 mans [Muñoz-Ortiz et al., 2024, Reinhart et al., 2025]. Just as different people can write in different  
163 styles, LLMs can also generate varied outputs. We think that what is commonly referred to as  
164 “LLM-generated text” is only a subset of the text that LLMs can potentially produce, and it’s likely  
165 the kind that corresponds to the most common and direct prompts. For instance, many detectors  
166 are trained on text generated by LLMs, which cannot represent all possibilities. Consequently, their  
167 detection capabilities are constrained. While different parameters can be set for various types of  
168 cases [Hans et al., 2024], such configurations can hardly cover all possible scenarios.

169 As we mentioned in related work, although some researchers have proposed benchmarks for detecting  
170 LLM-generated text [Liyanae et al., 2022, Wang et al., 2024b, He et al., 2024, Tao et al., 2024, Wu  
171 et al., 2024], their adoption has not yet become as widespread as other well-known LLM benchmarks,  
172 such as GLUE [Wang et al., 2018] and MMLU [Hendrycks et al., 2020]. Although these benchmarks  
173 have also faced some criticism [Hadi et al., 2023], there is still no highly universal benchmark for  
174 detecting LLM-generated texts.

175 Besides, due to the diversity of usage scenarios and the continuous updates of LLMs, a gold-standard  
176 benchmark is hard to realize, may even remain permanently absent.

## 177 5.2 Inherent Imperfection of Detection Tools

178 No existing LLM-detection system is infallible. In real-world conditions, a detection accuracy of  
179 85% is typically considered outstanding. Yet, this figure necessarily implies a 15% error rate, which  
180 may include both false positives and false negatives.

181 False positives—in which human-written content is incorrectly flagged as machine-generated—are  
182 particularly problematic in high-stakes contexts such as academic integrity investigations. And the  
183 problem of false positives has already been observed and discussed. For example, Rashidi et al.  
184 [2023] found that the AI text detector erroneously identified up to 8% of the known real abstracts as  
185 AI-generated text, and Giray [2024] states that false positives disproportionately affect non-native  
186 English speakers and scholars with distinctive writing styles. In addition, Tufts et al. [2024] think that  
187 adversarial attacks can easily bypass these detectors, and balancing high sensitivity with a reasonable  
188 true positive rate remains challenging.

189 Accusing someone of misconduct based on an imperfect tool can lead to unjust outcomes, reputational  
190 damage, and institutional distrust. Therefore, even detectors with relatively high accuracy present  
191 significant risks when used for evaluative or disciplinary purposes.

192 Experiments also show that certain detectors may exhibit bias against non-native English writ-  
193 ers [Liang et al., 2023] or against certain demographic groups [Kadoma et al., 2025]. The analyses  
194 from Li and Wan [2025] revealed that all the detectors they tested are highly sensitive to CEFR level  
195 and language environment. With LLMs being so widely used in academia [Eger et al., 2025, Russell  
196 et al., 2025], detecting AI-generated text must be handled with extreme care.

197 These detectors also face numerous other challenges, including difficulties in short-text detec-  
198 tion [Gameiro et al., 2024, Shi et al., 2024] and the issues of modification and classification that we  
199 will discuss later. As such, current detectors are far from perfect and may never achieve perfection in  
200 the future either.

### 201 5.3 Poor Robustness to Textual Modifications

202 There have always been many doubts about the effectiveness of these detectors [Sadasivan et al.,  
203 2023, Weber-Wulff et al., 2023]. Another issue pertains to the brittleness of these tools in realistic  
204 scenarios. An early study has shown that while humans can reliably detect poetry produced by GPT-2,  
205 but they struggle to accurately recognize it after human selection [Köbis and Mossink, 2021]. If  
206 post-generation modifications are taken into account, the detection process should become more  
207 challenging.

208 Most current detectors are trained to recognize text that has been directly generated by an LLM  
209 without post-editing. While some recent systems claim to maintain performance when the LLM-  
210 generated text is lightly modified, empirical evidence shows that detection accuracy tends to decline  
211 as the extent of human revision increases.

212 In practice, LLM-generated content is often edited, paraphrased, or interwoven with human-written  
213 material, especially in academic contexts. Consequently, the tools' applicability to real-world use  
214 cases remains limited. This limitation exacerbates the concerns raised in the first point, as reliance  
215 on imperfect systems in nuanced or ambiguous situations increases the likelihood of erroneous  
216 judgments.

### 217 5.4 A Wide Variety of Use Cases and the Limits of Binary Classification

218 A fourth, and perhaps more fundamental, concern lies in the heterogeneity of LLM-generated content.

219 The ethical implications of LLM use depend heavily on the context and intent of usage. For instance, a  
220 researcher who uses LLMs to generate entire manuscripts with minimal intellectual input contributes  
221 to the proliferation of unoriginal work, thereby burdening peer-review systems and undermining the  
222 credibility of scholarly communication.

223 Such practices are clearly unethical. In contrast, a non-native speaker might use an LLM to translate,  
224 rephrase, or refine a manuscript that is otherwise the product of original research. In this case, the  
225 LLM acts as a language aid rather than a generator of substantive content. Yet most detection systems  
226 treat these qualitatively different scenarios in the same manner, reducing the complex spectrum  
227 of authorship to a binary classification of "human-written" versus "machine-generated". Similar  
228 problems have also been noted in very recent studies [Lepp and Smith, 2025].

229 Generally, most studies focus on the binary classification problem of determining whether a given text  
230 is generated by LLMs. While some detection methods could achieve good results on given datasets,  
231 the scenario becomes more much complicated in real-world settings. For example, people could  
232 edit LLM-generated text or mix it with human written text, which has also attracted considerable  
233 attention [Zhang et al., 2024a, Abassy et al., 2024, Kumar et al., 2025]. Only a small number of  
234 researchers have tried to identify specific roles of LLM in content generation [Cheng et al., 2025],  
235 and no universally accepted approaches have been established.

### 236 5.5 Diversity in LLMs

237 Even without considering the usage scenarios noted before, different LLMs generate text in different  
238 styles [Rosenfeld and Lazebnik, 2024, Sun et al., 2025]. Empirical studies have consistently demon-  
239 strated that different LLMs exhibit distinct stylistic patterns fingerprints, which could even be used  
240 for classification [McGovern et al., 2024, Sun et al., 2025, Bitton et al., 2025].

241 Studies indicate that the detectability of texts depends on the LLM used for text generation [Antoun  
242 et al., 2023b]. For example, Wu et al. [2024] pointed out that the Binoculars [Hans et al., 2024] only  
243 achieved a 55.15% AUROC in detecting texts generated by Claude, while for texts generated by  
244 several other models, it reached at least 88%. A comparable point is reflected in Table 2.

245 Detectors may more easily flag text from older and smaller models [Elkhatat et al., 2023, Saha  
246 and Feizi, 2025]. As we all know, the development of LLMs has not stopped, so the timeliness of  
247 detectors is also another challenge. Obviously, the same LLM can produce different texts in response  
248 to different prompts for the same task, as we have shown before. **Although these detectors may still  
249 be applicable in certain scenarios, their use requires greater caution.**

## 250 5.6 Others

251 Those familiar with LLMs and detectors are aware of the potential issues, but the public tends to be  
252 easily drawn to these numbers and brief conclusion. The lack of detector interpretability represents  
253 another concern [Ji et al., 2024], severely limiting the ability to provide transparent explanations to  
254 the public.

255 The appropriate use of LLMs has now been widely accepted, such as in the NeurIPS submission  
256 process <sup>2</sup>. In addition to the examples given earlier, the traces of LLM-generated text have now been  
257 found in various fields, such as student essays’ answers [Leppänen et al., 2025] and words used in  
258 speaking [Yakura et al., 2024, Geng et al., 2024].

259 The abuse and misuse of these detectors can create ethical risks. Meanwhile, the numerical effective-  
260 ness of LLM-generated text detectors is declining. On the one hand, human may be influenced by  
261 LLMs and may create text resembling LLM-generated text. On the other hand, people may also  
262 adapt their language to bypass LLM detection tools [Geng and Trotta, 2025].

263 Therefore, **when interpreting the detection results of LLM-generated text, it is necessary to**  
264 **explicitly specify which kind of subset serves as the reference to establish the detector.**

## 265 6 Positive Impact of LLM Usage

266 The social impact of of LLMs has already been considered [Solaiman et al., 2019].

267 While LLM-generated text is frequently the subject of criticism—particularly due to concerns around  
268 academic dishonesty, plagiarism, and fraud, which have led to the development of various detection  
269 tools—it is equally important to emphasize the legitimate and ethical uses of large language models.  
270 As discussed earlier, LLMs can play a valuable role in numerous contexts. For instance, they help  
271 bridge linguistic divides by enabling non-native speakers to produce coherent and idiomatic texts in  
272 English or other target languages, thereby supporting greater inclusivity in academic and professional  
273 communication. They also facilitate high-quality machine translation, making content in multiple  
274 languages more accessible, and allow for the efficient synthesis of large textual corpora, which can  
275 aid research and knowledge production.

276 Rather than focusing solely on the detection and policing of LLM-generated text, it may be more  
277 productive to advocate for transparency regarding their use. In academic publishing, for example,  
278 it is increasingly common to disclose how LLMs have assisted in drafting, editing, or rephrasing  
279 portions of a manuscript.

280 Such uses are generally limited to improving expression or exploring alternative formulations; the  
281 substantive intellectual work remains the responsibility of human authors. Importantly, LLMs should  
282 not be considered co-authors nor used to autonomously generate scientific content in its entirety. Clear  
283 guidelines and disclosures can thus help normalize the ethical integration of LLMs into scholarly  
284 workflows without undermining academic integrity.

285 People began discussing ChatGPT’s positive impact shortly after its emergence [Kasneci et al., 2023].  
286 Non-native English speakers have to put in more effort as scientists, and there has been discrimination  
287 in the past [Amano et al., 2023, Lepp and Smith, 2025]. Automatic editing methods have shown  
288 promise in improving alignment between LLM-generated and human-written text [Chakrabarty et al.,  
289 2024]. If LLMs are applied properly and people assess detection tools reasonably, their positive  
290 influence can be greatly amplified.

## 291 7 Alternative Views

292 Researchers have not yet reached full agreement on the detectability of LLM-generated text.

293 Chakraborty et al. [2024] claim in their position paper: “Despite ongoing debate about the feasibility  
294 of such differentiation, we present evidence supporting its consistent achievability, except when  
295 human and machine text distributions are indistinguishable across their entire support. Drawing from

---

<sup>2</sup><https://neurips.cc/Conferences/2025/LLM>

296 information theory, we argue that as machine-generated text approximates human-like quality, the  
297 sample size needed for detection increases.”

298 But we believe that LLMs are fully capable of generating text that is nearly indistinguishable from  
299 human-written content. Furthermore, practical observations have shown that humans possess the  
300 capacity to identify LLM-generated text with reasonable accuracy [Russell et al., 2025], and such  
301 coevolution may already be occurring [Geng and Trotta, 2025]. These challenges in these real-world  
302 data cannot be resolved by increasing the sample size.

303 The key disagreement among researchers may not be technical in nature, but rather stems from differ-  
304 ing perspectives on human intervention. Take watermarking studies as an example, if people edit the  
305 generated text (which is simple to do), the watermark’s reliability may be greatly weakened [Dathathri  
306 et al., 2024].

307 There are also researchers who share similar views with us. For example, Zhang et al. [2024b] argue  
308 that "We believe that the issue of AI-generated text detection remains an unresolved challenge. As  
309 LLMs become increasingly powerful and humans become more proficient in using them, it is even  
310 less likely to detect AI text in the future." And Nicks et al. [2023] “advise against continued reliance  
311 on LLM-generated text detectors”.

## 312 **8 Future Perspectives and Predictions**

313 LLMs were compared to stochastic parrot [Bender et al., 2021] a couple of years ago, but their  
314 capabilities are gradually being recognized [Srivastava et al., 2022], and now their competencies have  
315 reached or even surpassed those of human experts in various fields. This leads us to speculate that  
316 future LLM-generated texts could surpass current versions in human resemblance. The persona effect  
317 could be a good example [Hu and Collier, 2024].

318 The gap between LLM-generated and human-written text is expected to narrow further. As noted  
319 by several researchers, humans can learn to detect AI-generated texts [Milička et al., 2025] and can  
320 become accurate and robust detector of LLM-generated text [Russell et al., 2025].

321 And therefore, we can find more cases of coevolution [Pedreschi et al., 2024, Geng and Trotta, 2025]  
322 between human and Human in the future. Consequently, in the future, detecting LLM-generated text  
323 may become less important, particularly in terms of numerical interpretation.

324 Accounting for model collapse [Shumailov et al., 2024, Guo et al., 2023] and knowledge collapse [Pe-  
325 terson, 2025], the detection results may become even more intriguing. Detection efforts should target  
326 substantive content (e.g., fact-checking) rather than linguistic characteristics [Schuster et al., 2020].

## 327 **9 Conclusions**

328 Given the current state-of-the-art, existing tools are not equipped to make such fine-grained dis-  
329 tinctions. They are structurally unable to assess the proportion, function, or ethical significance of  
330 LLM contributions in a given text. As a result, the development and deployment of LLM-generated  
331 text detectors raise serious concerns, not only due to technical limitations but also because they  
332 risk enforcing overly simplistic and potentially unjust frameworks for evaluating authorship and  
333 intellectual responsibility.

334 And some of these difficulties are simply unavoidable, not merely temporary technical challenges.  
335 Texts generated by LLMs and those generated by humans often overlap greatly, with no obvious  
336 separation. Moreover, as LLMs become more widely used and people may learn from their outputs,  
337 the difference between them may get further smaller. Since text is different from images, it is difficult  
338 to balance both the watermark and the original textual information.

339 While acknowledging that detectors for LLM-generated text can serve a purpose in certain scenarios,  
340 we recommend using them with greater caution. And the detection results should only be used as a  
341 reference rather than the decisive indicator.



## 342 References

- 343 Mervat Abassy, Kareem Elozeiri, Alexander Aziz, Minh Ngoc Ta, Raj Vardhan Tomar, Bimar-  
344 sha Adhikari, Saad El Dine Ahmed, Yuxia Wang, Osama Mohammed Afzal, Zhuohan Xie,  
345 et al. Llm-detectaive: a tool for fine-grained machine-generated text detection. *arXiv preprint*  
346 *arXiv:2408.04284*, 2024.
- 347 Sara Abdali, Richard Anarfi, CJ Barberan, and Jia He. Decoding the ai pen: Techniques and  
348 challenges in detecting ai-generated text. In *Proceedings of the 30th ACM SIGKDD Conference on*  
349 *Knowledge Discovery and Data Mining*, pages 6428–6436, 2024.
- 350 Tatsuya Amano, Valeria Ramírez-Castañeda, Violeta Berdejo-Espinola, Israel Borokini, Shawan  
351 Chowdhury, Marina Golivets, Juan David González-Trujillo, Flavia Montaña-Centellas, Kumar  
352 Paudel, Rachel Louise White, et al. The manifold costs of being a non-native english speaker in  
353 science. *PLoS Biology*, 21(7):e3002184, 2023.
- 354 Wissam Antoun, Virginie Mouilleron, Benoît Sagot, and Djamé Seddah. Towards a robust detection  
355 of language model generated text: is chatgpt that easy to detect? *arXiv preprint arXiv:2306.05871*,  
356 2023a.
- 357 Wissam Antoun, Benoît Sagot, and Djamé Seddah. From text to source: Results in detecting large  
358 language model-generated content. *arXiv preprint arXiv:2309.13322*, 2023b.
- 359 Guangsheng Bao, Yanbin Zhao, Zhiyang Teng, Linyi Yang, and Yue Zhang. Fast-detectgpt: Efficient  
360 zero-shot detection of machine-generated text via conditional probability curvature. *arXiv preprint*  
361 *arXiv:2310.05130*, 2023.
- 362 Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the  
363 dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM*  
364 *conference on fairness, accountability, and transparency*, pages 610–623, 2021.
- 365 Amrita Bhattacharjee, Tharindu Kumarage, Raha Moraffah, and Huan Liu. Conda: Contrastive  
366 domain adaptation for ai-generated text detection. *arXiv preprint arXiv:2309.03992*, 2023.
- 367 Yehonatan Bitton, Elad Bitton, and Shai Nisan. Detecting stylistic fingerprints of large language  
368 models. *arXiv preprint arXiv:2503.01659*, 2025.
- 369 Creston Brooks, Samuel Eggert, and Denis Peskoff. The rise of ai-generated content in wikipedia.  
370 *arXiv preprint arXiv:2410.08044*, 2024.
- 371 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
372 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
373 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 374 Shuyang Cai and Wanyun Cui. Evade chatgpt detectors via a single space. *arXiv preprint*  
375 *arXiv:2307.02599*, 2023.
- 376 Tuhin Chakrabarty, Philippe Laban, and Chien-Sheng Wu. Can ai writing be salvaged? mitigating  
377 idiosyncrasies and improving human-ai alignment in the writing process through edits. *arXiv*  
378 *preprint arXiv:2409.14509*, 2024.
- 379 Souradip Chakraborty, Amrit Bedi, Sicheng Zhu, Bang An, Dinesh Manocha, and Furong Huang.  
380 Position: On the possibilities of ai-generated text detection. In *Forty-first International Conference*  
381 *on Machine Learning*, 2024.
- 382 Zihao Cheng, Li Zhou, Feng Jiang, Benyou Wang, and Haizhou Li. Beyond binary: Towards  
383 fine-grained llm-generated text detection via role recognition and involvement measurement. In  
384 *Proceedings of the ACM on Web Conference 2025*, pages 2677–2688, 2025.
- 385 Miranda Christ, Sam Gunn, and Or Zamir. Undetectable watermarks for language models. In *The*  
386 *Thirty Seventh Annual Conference on Learning Theory*, pages 1125–1139. PMLR, 2024.
- 387 Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A Smith.  
388 All that’s human’s not gold: Evaluating human evaluation of generated text. *arXiv preprint*  
389 *arXiv:2107.00061*, 2021.

- 390 Joseph Cornelius, Oscar Lithgow-Serrano, Sandra Mitrović, Ljiljana Dolamic, and Fabio Rinaldi.  
391 Bust: Benchmark for the evaluation of detectors of llm-generated text. In *Proceedings of the*  
392 *2024 Conference of the North American Chapter of the Association for Computational Linguistics:*  
393 *Human Language Technologies (Volume 1: Long Papers)*, pages 8029–8057, 2024.
- 394 Evan N Crothers, Nathalie Japkowicz, and Herna L Viktor. Machine-generated text: A comprehensive  
395 survey of threat models and detection methods. *IEEE Access*, 11:70977–71002, 2023.
- 396 Sumanth Dathathri, Abigail See, Sumedh Ghaisas, Po-Sen Huang, Rob McAdam, Johannes Welbl,  
397 Vandana Bachani, Alex Kaskasoli, Robert Stanforth, Tatiana Matejovicova, et al. Scalable water-  
398 marking for identifying large language model outputs. *Nature*, 634(8035):818–823, 2024.
- 399 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
400 bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of*  
401 *the North American chapter of the association for computational linguistics: human language*  
402 *technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.
- 403 Jad Doughman, Osama Mohammed Afzal, Hawau Olamide Toyin, Shady Shehata, Preslav Nakov,  
404 and Zeerak Talat. Exploring the limitations of detecting machine-generated text. *arXiv preprint*  
405 *arXiv:2406.11073*, 2024.
- 406 Liam Dugan, Daphne Ippolito, Arun Kirubarajan, Sherry Shi, and Chris Callison-Burch. Real or  
407 fake text?: Investigating human ability to detect boundaries between human-written and machine-  
408 generated text. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages  
409 12763–12771, 2023.
- 410 Steffen Eger, Yong Cao, Jennifer D’Souza, Andreas Geiger, Christian Greisinger, Stephanie Gross,  
411 Yufang Hou, Brigitte Krenn, Anne Lauscher, Yizhi Li, et al. Transforming science with large  
412 language models: A survey on ai-assisted scientific discovery, experimentation, content generation,  
413 and evaluation. *arXiv preprint arXiv:2502.05151*, 2025.
- 414 Ahmed M Elkhataf, Khaled Elsaid, and Saeed Almeer. Evaluating the efficacy of ai content detection  
415 tools in differentiating between human and ai-generated text. *International Journal for Educational*  
416 *Integrity*, 19(1):17, 2023.
- 417 Kathleen C Fraser, Hillary Dawkins, and Svetlana Kiritchenko. Detecting ai-generated text: Factors  
418 influencing detectability with current methods. *arXiv preprint arXiv:2406.15583*, 2024.
- 419 Margherita Gambini, Tiziano Fagni, Fabrizio Falchi, and Maurizio Tesconi. On pushing deepfake  
420 tweet detection capabilities to the limits. In *Proceedings of the 14th ACM Web Science Conference*  
421 *2022*, pages 154–163, 2022.
- 422 Henrique Da Silva Gameiro, Andrei Kucharavy, and Ljiljana Dolamic. Llm detectors still fall short of  
423 real world: Case of llm-generated short news-like posts. *arXiv preprint arXiv:2409.03291*, 2024.
- 424 Sebastian Gehrmann, Hendrik Strobelt, and Alexander M Rush. Gltr: Statistical detection and  
425 visualization of generated text. *arXiv preprint arXiv:1906.04043*, 2019.
- 426 Mingmeng Geng and Roberto Trotta. Is chatgpt transforming academics’ writing style? *arXiv*  
427 *preprint arXiv:2404.08627*, 2024.
- 428 Mingmeng Geng and Roberto Trotta. Human-llm coevolution: Evidence from academic writing.  
429 *arXiv preprint arXiv:2502.09606*, 2025.
- 430 Mingmeng Geng, Caixi Chen, Yanru Wu, Dongping Chen, Yao Wan, and Pan Zhou. The impact of  
431 large language models in academia: from writing to speaking. *arXiv preprint arXiv:2409.13686*,  
432 2024.
- 433 Louie Giray. The problem with false positives: Ai detection unfairly accuses scholars of ai plagiarism.  
434 *The Serials Librarian*, 85(5-6):181–189, 2024.
- 435 Alexei Grinbaum and Laurynas Adomaitis. The ethical need for watermarks in machine-generated  
436 language. *arXiv preprint arXiv:2209.03118*, 2022.

- 437 German Gritsai, Anastasia Voznyuk, Andrey Grabovoy, and Yury Chekhovich. Are ai detectors  
438 good enough? a survey on quality of datasets with machine-generated texts. *arXiv preprint*  
439 *arXiv:2410.14677*, 2024.
- 440 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
441 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms  
442 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 443 Yanzhu Guo, Guokan Shang, Michalis Vazirgiannis, and Chloé Clavel. The curious decline of  
444 linguistic diversity: Training language models on synthetic text. *arXiv preprint arXiv:2311.09807*,  
445 2023.
- 446 Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muham-  
447 mad Bilal Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili, et al. A survey on large language  
448 models: Applications, challenges, limitations, and practical usage. *Authorea Preprints*, 2023.
- 449 Abhimanyu Hans, Avi Schwarzschild, Valeriia Cherepanova, Hamid Kazemi, Aniruddha Saha, Micah  
450 Goldblum, Jonas Geiping, and Tom Goldstein. Spotting llms with binoculars: Zero-shot detection  
451 of machine-generated text. *arXiv preprint arXiv:2401.12070*, 2024.
- 452 Xinlei He, Xinyue Shen, Zeyuan Chen, Michael Backes, and Yang Zhang. Mgtbench: Benchmarking  
453 machine-generated text detection. In *Proceedings of the 2024 on ACM SIGSAC Conference on*  
454 *Computer and Communications Security*, pages 2251–2265, 2024.
- 455 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
456 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*  
457 *arXiv:2009.03300*, 2020.
- 458 Tiancheng Hu and Nigel Collier. Quantifying the persona effect in llm simulations. *arXiv preprint*  
459 *arXiv:2402.10811*, 2024.
- 460 Siming Huang, Yuliang Xu, Mingmeng Geng, Yao Wan, and Dongping Chen. Wikipedia in the era  
461 of llms: Evolution and risks. *arXiv preprint arXiv:2503.02879*, 2025.
- 462 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-  
463 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint*  
464 *arXiv:2410.21276*, 2024.
- 465 Daphne Ippolito, Daniel Duckworth, Chris Callison-Burch, and Douglas Eck. Automatic detection of  
466 generated text is easiest when humans are fooled. *arXiv preprint arXiv:1911.00650*, 2019.
- 467 Jiazhou Ji, Ruizhe Li, Shujun Li, Jie Guo, Weidong Qiu, Zheng Huang, Chiyu Chen, Xiaoyu Jiang,  
468 and Xinru Lu. Detecting machine-generated texts: Not just "ai vs humans" and explainability is  
469 complicated. *arXiv preprint arXiv:2406.18259*, 2024.
- 470 Nikola Jovanović, Robin Staab, and Martin Vechev. Watermark stealing in large language models.  
471 *arXiv preprint arXiv:2402.19361*, 2024.
- 472 Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert  
473 McHardy. Challenges and applications of large language models. *arXiv preprint arXiv:2307.10169*,  
474 2023.
- 475 Kowe Kadoma, Danaë Metaxa, and Mor Naaman. Generative ai and perceptual harms: Who's  
476 suspected of using llms? In *Proceedings of the 2025 CHI Conference on Human Factors in*  
477 *Computing Systems*, pages 1–17, 2025.
- 478 Sujita Kumar Kar, Teena Bansal, Sumit Modi, and Amit Singh. How sensitive are the free ai-detector  
479 tools in detecting ai-generated texts? a comparison of popular ai-detector tools. *Indian Journal of*  
480 *Psychological Medicine*, page 02537176241247934, 2024.
- 481 Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank  
482 Fischer, Urs Gasser, Georg Groh, Stephan Günemann, Eyke Hüllermeier, et al. Chatgpt for good?  
483 on opportunities and challenges of large language models for education. *Learning and individual*  
484 *differences*, 103:102274, 2023.

- 485 John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A  
486 watermark for large language models. In *International Conference on Machine Learning*, pages  
487 17061–17084. PMLR, 2023.
- 488 Nils Köbis and Luca D Mossink. Artificial intelligence versus maya angelou: Experimental evidence  
489 that people cannot differentiate ai-generated from human-written poetry. *Computers in human*  
490 *behavior*, 114:106553, 2021.
- 491 Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. Paraphrasing  
492 evades detectors of ai-generated text, but retrieval is an effective defense. *Advances in Neural*  
493 *Information Processing Systems*, 36:27469–27500, 2023.
- 494 Sandeep Kumar, Samarth Garg, Sagnik Sengupta, Tirthankar Ghosal, and Asif Ekbal. Mixrevdetect:  
495 Towards detecting ai-generated content in hybrid peer reviews. In *Proceedings of the 2025 Con-*  
496 *ference of the Nations of the Americas Chapter of the Association for Computational Linguistics:*  
497 *Human Language Technologies (Volume 2: Short Papers)*, pages 944–953, 2025.
- 498 Tharindu Kumarage, Joshua Garland, Amrita Bhattacharjee, Kirill Trapeznikov, Scott Ruston,  
499 and Huan Liu. Stylometric detection of ai-generated text in twitter timelines. *arXiv preprint*  
500 *arXiv:2303.03697*, 2023.
- 501 Tharindu Kumarage, Garima Agrawal, Paras Sheth, Raha Moraffah, Aman Chadha, Joshua Gar-  
502 land, and Huan Liu. A survey of ai-generated text forensic systems: Detection, attribution, and  
503 characterization. *arXiv preprint arXiv:2403.01152*, 2024.
- 504 Hiu Ting Lau and Arkaitz Zubiaga. Understanding the effects of human-written paraphrases in  
505 llm-generated text detection. *arXiv preprint arXiv:2411.03806*, 2024.
- 506 Haley Lepp and Daniel Scott Smith. " you cannot sound like gpt": Signs of language discrimination  
507 and resistance in computer science publishing. *arXiv preprint arXiv:2505.08127*, 2025.
- 508 Leo Leppänen, Lili Aunimo, Arto Hellas, Jukka K Nurminen, and Linda Mannila. How large  
509 language models are changing mooc essay answers: A comparison of pre-and post-llm responses.  
510 *arXiv preprint arXiv:2504.13038*, 2025.
- 511 Jiatao Li and Xiaojun Wan. Who writes what: Unveiling the impact of author roles on ai-generated  
512 text detection. *arXiv preprint arXiv:2502.12611*, 2025.
- 513 Weixin Liang, Mert Yuksekgonul, Yining Mao, Eric Wu, and James Zou. Gpt detectors are biased  
514 against non-native english writers. *Patterns*, 4(7), 2023.
- 515 Weixin Liang, Zachary Izzo, Yaohui Zhang, Haley Lepp, Hancheng Cao, Xuandong Zhao, Lingjiao  
516 Chen, Haotian Ye, Sheng Liu, Zhi Huang, et al. Monitoring ai-modified content at scale: A case  
517 study on the impact of chatgpt on ai conference peer reviews. *arXiv preprint arXiv:2403.07183*,  
518 2024.
- 519 Weixin Liang, Yaohui Zhang, Mihai Codreanu, Jiayu Wang, Hancheng Cao, and James Zou. The  
520 widespread adoption of large language model-assisted writing across society. *arXiv preprint*  
521 *arXiv:2502.09747*, 2025.
- 522 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,  
523 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*  
524 *arXiv:2412.19437*, 2024.
- 525 Zhengliang Liu, Yue Huang, Xiaowei Yu, Lu Zhang, Zihao Wu, Chao Cao, Haixing Dai, Lin Zhao,  
526 Yiwei Li, Peng Shu, et al. Deid-gpt: Zero-shot medical text de-identification by gpt-4. *arXiv*  
527 *preprint arXiv:2303.11032*, 2023.
- 528 Vijini Liyanage, Davide Buscaldi, and Adeline Nazarenko. A benchmark corpus for the detection of  
529 automatically generated text in academic publications. *arXiv preprint arXiv:2202.02013*, 2022.
- 530 Shixuan Ma and Quan Wang. Zero-shot detection of llm-generated text using token cohesiveness.  
531 *arXiv preprint arXiv:2409.16914*, 2024.

- 532 Dominik Macko, Robert Moro, Adaku Uchendu, Ivan Srba, Jason Samuel Lucas, Michiharu Ya-  
533 mashita, Nafis Irtiza Tripto, Dongwon Lee, Jakub Simko, and Maria Bielikova. Authorship  
534 obfuscation in multilingual machine-generated text detection. *arXiv preprint arXiv:2401.07867*,  
535 2024.
- 536 Hope McGovern, Rickard Stureborg, Yoshi Suhara, and Dimitris Alikaniotis. Your large language  
537 models are leaving fingerprints. *arXiv preprint arXiv:2405.14057*, 2024.
- 538 Jiří Milička, Anna Marklová, Ondřej Drobil, and Eva Pospíšilová. Humans can learn to detect  
539 ai-generated texts, or at least learn when they can't. *arXiv preprint arXiv:2505.01877*, 2025.
- 540 Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D Manning, and Chelsea Finn.  
541 Detectgpt: Zero-shot machine-generated text detection using probability curvature. In *International*  
542 *Conference on Machine Learning*, pages 24950–24962. PMLR, 2023.
- 543 Alberto Muñoz-Ortiz, Carlos Gómez-Rodríguez, and David Vilares. Contrasting linguistic patterns  
544 in human and llm-generated news text. *Artificial Intelligence Review*, 57(10):265, 2024.
- 545 Charlotte Nicks, Eric Mitchell, Rafael Rafailov, Archit Sharma, Christopher D Manning, Chelsea  
546 Finn, and Stefano Ermon. Language model detectors are easily optimized against. In *The twelfth*  
547 *international conference on learning representations*, 2023.
- 548 Dino Pedreschi, Luca Pappalardo, Emanuele Ferragina, Ricardo Baeza-Yates, Albert-László Barabási,  
549 Frank Dignum, Virginia Dignum, Tina Eliassi-Rad, Fosca Giannotti, János Kertész, et al. Human-ai  
550 coevolution. *Artificial Intelligence*, page 104244, 2024.
- 551 Andrew J Peterson. Ai and the problem of knowledge collapse. *AI & SOCIETY*, pages 1–21, 2025.
- 552 Shushanta Pudasaini, Luis Miralles-Pechuán, David Lillis, and Marisa Llorens Salvador. Survey on  
553 plagiarism detection in large language models: The impact of chatgpt and gemini on academic  
554 integrity. *arXiv preprint arXiv:2407.13105*, 2024.
- 555 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
556 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 557 Hooman H Rashidi, Brandon D Fennell, Samer Albahra, Bo Hu, and Tom Gorbett. The chatgpt  
558 conundrum: Human-generated scientific manuscripts misidentified as ai creations by ai text  
559 detection tool. *Journal of Pathology Informatics*, 14:100342, 2023.
- 560 Alex Reinhart, Ben Markey, Michael Laudenbach, Kachata Pantusen, Ronald Yurko, Gordon  
561 Weinberg, and David West Brown. Do llms write like humans? variation in grammatical and  
562 rhetorical styles. *Proceedings of the National Academy of Sciences*, 122(8):e2422455122, 2025.
- 563 Ariel Rosenfeld and Teddy Lazebnik. Whose llm is it anyway? linguistic comparison and llm  
564 attribution for gpt-3.5, gpt-4 and bard. *arXiv preprint arXiv:2402.14533*, 2024.
- 565 Jenna Russell, Marzena Karpinska, and Mohit Iyyer. People who frequently use chatgpt for writing  
566 tasks are accurate and robust detectors of ai-generated text. *arXiv preprint arXiv:2501.15654*,  
567 2025.
- 568 Vinu Sankar Sadasivan, Aounon Kumar, Sriram Balasubramanian, Wenxiao Wang, and Soheil Feizi.  
569 Can ai-generated text be reliably detected? *arXiv preprint arXiv:2303.11156*, 2023.
- 570 Shoumik Saha and Soheil Feizi. Almost ai, almost human: The challenge of detecting ai-polished  
571 writing. *arXiv preprint arXiv:2502.15666*, 2025.
- 572 Tal Schuster, Roei Schuster, Darsh J Shah, and Regina Barzilay. The limitations of stylometry for  
573 detecting machine-generated fake news. *Computational Linguistics*, 46(2):499–510, 2020.
- 574 Yuhui Shi, Qiang Sheng, Juan Cao, Hao Mi, Beizhe Hu, and Danding Wang. Ten words only still  
575 help: Improving black-box ai-generated text detection via proxy-guided efficient re-sampling.  
576 *arXiv preprint arXiv:2402.09199*, 2024.
- 577 Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal. Ai  
578 models collapse when trained on recursively generated data. *Nature*, 631(8022):755–759, 2024.

- 579 Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askill, Ariel Herbert-Voss, Jeff Wu, Alec  
580 Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, et al. Release strategies and the social  
581 impacts of language models. *arXiv preprint arXiv:1908.09203*, 2019.
- 582 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam  
583 Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the  
584 imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint*  
585 *arXiv:2206.04615*, 2022.
- 586 Jinyan Su, Terry Yue Zhuo, Di Wang, and Preslav Nakov. Detectllm: Leveraging log rank information  
587 for zero-shot detection of machine-generated text. *arXiv preprint arXiv:2306.05540*, 2023.
- 588 Mingjie Sun, Yida Yin, Zhiqiu Xu, J Zico Kolter, and Zhuang Liu. Idiosyncrasies in large language  
589 models. *arXiv preprint arXiv:2502.12150*, 2025.
- 590 Zhen Sun, Zongmin Zhang, Xinyue Shen, Ziyi Zhang, Yule Liu, Michael Backes, Yang Zhang, and  
591 Xinlei He. Are we in the ai-generated text world already? quantifying and monitoring aigt on  
592 social media. *arXiv preprint arXiv:2412.18148*, 2024.
- 593 Ruixiang Tang, Yu-Neng Chuang, and Xia Hu. The science of detecting llm-generated text. *Commu-*  
594 *nications of the ACM*, 67(4):50–59, 2024.
- 595 Zhen Tao, Zhiyu Li, Dinghao Xi, and Wei Xu. Cudrt: Benchmarking the detection of human vs. large  
596 language models generated texts. *arXiv preprint arXiv:2406.09056*, 2024.
- 597 Brian Tufts, Xuandong Zhao, and Lei Li. A practical examination of ai-generated text detectors for  
598 large language models. *arXiv preprint arXiv:2412.05139*, 2024.
- 599 Jan Philip Wahle, Terry Ruas, Frederic Kirstein, and Bela Gipp. How large language models are  
600 transforming machine-paraphrased plagiarism. *arXiv preprint arXiv:2210.03568*, 2022.
- 601 William H Walters. The effectiveness of software designed to detect ai-generated writing: A  
602 comparison of 16 ai text detectors. *Open Information Science*, 7(1):20220158, 2023.
- 603 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue:  
604 A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint*  
605 *arXiv:1804.07461*, 2018.
- 606 Quan Wang and Haoran Li. On continually tracing origins of llm-generated text and its application in  
607 detecting cheating in student coursework. *Big Data and Cognitive Computing*, 9(3):50, 2025.
- 608 Rongsheng Wang, Haoming Chen, Ruizhe Zhou, Han Ma, Yaofei Duan, Yanlan Kang, Songhua  
609 Yang, Baoyu Fan, and Tao Tan. Llm-detector: Improving ai-generated chinese text detection with  
610 open-source llm instruction tuning. *arXiv preprint arXiv:2402.01158*, 2024a.
- 611 Yuxia Wang, Jonibek Mansurov, Petar Ivanov, Jinyan Su, Artem Shelmanov, Akim Tsvigun,  
612 Osama Mohammed Afzal, Tarek Mahmoud, Giovanni Puccetti, Thomas Arnold, et al. M4gt-  
613 bench: Evaluation benchmark for black-box machine-generated text detection. *arXiv preprint*  
614 *arXiv:2402.11175*, 2024b.
- 615 Yuxia Wang, Artem Shelmanov, Jonibek Mansurov, Akim Tsvigun, Vladislav Mikhailov, Rui Xing,  
616 Zhuohan Xie, Jiahui Geng, Giovanni Puccetti, Ekaterina Artemova, et al. Genai content detection  
617 task 1: English and multilingual machine-generated text detection: Ai vs. human. *arXiv preprint*  
618 *arXiv:2501.11012*, 2025.
- 619 Debora Weber-Wulff, Alla Anohina-Naumeca, Sonja Bjelobaba, Tomáš Foltýnek, Jean Guerrero-Dib,  
620 Olumide Popoola, Petr Šigut, and Lorna Waddington. Testing of detection tools for ai-generated  
621 text. *International Journal for Educational Integrity*, 19(1):1–39, 2023.
- 622 Junchao Wu, Runzhe Zhan, Derek Wong, Shu Yang, Xinyi Yang, Yulin Yuan, and Lidia Chao.  
623 Detectrl: Benchmarking llm-generated text detection in real-world scenarios. *Advances in Neural*  
624 *Information Processing Systems*, 37:100369–100401, 2024.

- 625 Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Lidia Sam Chao, and Derek Fai Wong. A  
626 survey on llm-generated text detection: Necessity, methods, and future directions. *Computational*  
627 *Linguistics*, pages 1–66, 2025.
- 628 Kangxi Wu, Liang Pang, Huawei Shen, Xueqi Cheng, and Tat-Seng Chua. Llmdet: A third party  
629 large language models generated text detection tool. *arXiv preprint arXiv:2305.15004*, 2023.
- 630 Hiromu Yakura, Ezequiel Lopez-Lopez, Levin Brinkmann, Ignacio Serna, Prateek Gupta, and Iyad  
631 Rahwan. Empirical evidence of large language model’s influence on human spoken communication.  
632 *arXiv preprint arXiv:2409.01754*, 2024.
- 633 Xianjun Yang, Liangming Pan, Xuandong Zhao, Haifeng Chen, Linda Petzold, William Yang Wang,  
634 and Wei Cheng. A survey on detection of llms-generated content. *arXiv preprint arXiv:2310.15654*,  
635 2023.
- 636 Xiao Yu, Kejiang Chen, Qi Yang, Weiming Zhang, and Nenghai Yu. Text fluoroscopy: Detecting  
637 llm-generated text through intrinsic features. In *Proceedings of the 2024 Conference on Empirical*  
638 *Methods in Natural Language Processing*, pages 15838–15846, 2024a.
- 639 Xiao Yu, Yuang Qi, Kejiang Chen, Guoqiang Chen, Xi Yang, Pengyuan Zhu, Xiuwei Shang, Weiming  
640 Zhang, and Nenghai Yu. Dpic: Decoupling prompt and intrinsic characteristics for llm generated  
641 text detection. *Advances in Neural Information Processing Systems*, 37:16194–16212, 2024b.
- 642 Wataru Zaitzu and Mingzhe Jin. Distinguishing chatgpt (-3.5,-4)-generated and human-written papers  
643 through japanese stylometric analysis. *PLoS One*, 18(8):e0288453, 2023.
- 644 Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and  
645 Yejin Choi. Defending against neural fake news. *Advances in neural information processing*  
646 *systems*, 32, 2019.
- 647 Qihui Zhang, Chujie Gao, Dongping Chen, Yue Huang, Yixin Huang, Zhenyang Sun, Shilin Zhang,  
648 Weiye Li, Zhengyan Fu, Yao Wan, et al. Llm-as-a-coauthor: Can mixed human-written and  
649 machine-generated text be detected? *arXiv preprint arXiv:2401.05952*, 2024a.
- 650 Yuehan Zhang, Yongqiang Ma, Jiawei Liu, Xiaozhong Liu, Xiaofeng Wang, and Wei Lu. Detection  
651 vs. anti-detection: Is text generated by ai detectable? In *International Conference on Information*,  
652 pages 209–222. Springer, 2024b.

653 **A Simulation Results**

- 654 • **DeepSeek-V3 (Polish the following passage:)**. Founded in 1987, the conference has  
655 grown into a premier multi-track interdisciplinary annual event. It features invited talks,  
656 demonstrations, symposia, and peer-reviewed oral and poster presentations. Alongside the  
657 conference, attendees can explore a professional exposition highlighting practical applica-  
658 tions of machine learning, a series of tutorials, and topical workshops that foster informal  
659 idea exchange.
- 660 • **DeepSeek-R1 (Polish the following passage:)**. Established in 1987, the conference has  
661 grown into a premier interdisciplinary event held annually, featuring multiple tracks of  
662 programming. Attendees engage with invited talks, demonstrations, symposia, and peer-  
663 reviewed research presented through both oral and poster sessions. Complementing the  
664 core program are a professional exposition showcasing real-world applications of machine  
665 learning, hands-on tutorials, and topical workshops designed to foster informal dialogue and  
666 collaborative idea exchange.
- 667 • **GPT-3.5 (Polish the following passage:)**. Established in 1987, the conference has evolved  
668 into a multi-track interdisciplinary annual event featuring invited talks, demonstrations,  
669 symposia, and oral and poster presentations of refereed papers. In addition, there is a  
670 professional exposition dedicated to machine learning in practice, a variety of tutorials, and  
671 topical workshops that offer a more relaxed environment for the sharing of ideas.
- 672 • **GPT-4o-mini (Polish the following passage:)**. Founded in 1987, the conference has evolved  
673 into a dynamic multi-track interdisciplinary annual event featuring invited talks, demonstra-  
674 tions, symposia, and both oral and poster presentations of refereed papers. Complementing  
675 the conference is a professional exposition dedicated to practical applications of machine  
676 learning, as well as a series of tutorials and topical workshops that create a more informal  
677 atmosphere for exchanging ideas.
- 678 • **GPT-4o (Polish the following passage:)**. Established in 1987, the conference has evolved  
679 into a comprehensive, multi-track interdisciplinary annual event. It features invited talks,  
680 demonstrations, symposia, and both oral and poster presentations of peer-reviewed papers.  
681 In addition to the main conference, there is a professional exposition dedicated to practical  
682 applications of machine learning, as well as a series of tutorials and topical workshops.  
683 These workshops offer a more informal setting for the exchange of ideas.
- 684 • **GPT-3.5 (Rewrite the following passage without using complex expressions:)**. The  
685 conference started in 1987 and has grown into a yearly event that covers different subjects  
686 and includes talks, demonstrations, symposiums, and presentations of reviewed papers.  
687 There is also an exhibition on machine learning, tutorials, and workshops for sharing ideas  
688 in a casual setting.
- 689 • **GPT-4o-mini (Rewrite the following passage without using complex expressions:)**. The  
690 conference started in 1987 and is now held every year. It features various sessions that  
691 include invited talks, demonstrations, and both oral and poster presentations of reviewed  
692 papers. There is also a professional exhibit that highlights real-world applications of machine  
693 learning, along with a series of tutorials and workshops that offer a more relaxed atmosphere  
694 for sharing ideas.
- 695 • **GPT-4o (Rewrite the following passage without using complex expressions:)**. The  
696 conference started in 1987 and now happens every year with multiple sessions. It features  
697 invited talks, demonstrations, discussions, and presentations of selected papers. There is  
698 also a professional expo about using machine learning, a series of tutorials, and workshops  
699 that offer a more relaxed space for sharing ideas.