

# Enhancing Robustness of LLM-Synthetic Text Detectors for Academic Writing: A Comprehensive Analysis

Anonymous ACL submission

## Abstract

The emergence of large language models (LLMs), such as Generative Pre-trained Transformer 4 (GPT-4) used by ChatGPT, has profoundly impacted the academic and broader community. While these models offer numerous advantages in revolutionizing work and study methods, they have also garnered significant attention due to their potential negative consequences. One example is generating academic reports or papers without or with a limited human contribution. Consequently, researchers have focused on developing detectors to address the misuse of LLMs. However, most existing works prioritize achieving higher accuracy on restricted datasets, neglecting the crucial aspect of generalizability. This limitation hinders their practical application in real-life scenarios where reliability is paramount. In this paper, we present a comprehensive analysis of the influence of prompts on the text generated by LLMs and highlight the potential lack of robustness in one of the current state-of-the-art GPT detectors. To mitigate these issues concerning the misuse of LLMs in academic writing, we propose a reference-based Siamese detector taking a pair of texts: one as the inquiry and the other as the reference. Our method effectively addresses the lack of robustness and significantly improves the baseline performances in challenging scenarios, increasing them by approximately 25% to 67%.

## 1 Introduction

Recently, the applications of large-scale language models, such as Open AI's GPT-4 (OpenAI, 2023), or Google's Pathways Language Model 2 (Anil et al., 2023), have become an integral part of people's lives and works, often being utilized unconsciously. From casual conversations with chatbots to accurately expressing search queries on search engines and relying on models like ChatGPT for writing assistants, LLMs have gained widespread usage due to their powerful performance. This

extensive application potential has attracted numerous companies to leverage LLMs for optimizing their services. However, while LLMs greatly facilitate daily activities, they pose significant security risks if maliciously exploited for attacks or deceptions. Consequently, with the growing popularity of LLMs, the importance of AI security has come to the forefront of people's attention.

Among the various security concerns, academic cheating stands out as a particularly grave issue. Within academia, universities face the most severe challenges in this regard. University students possess the necessary expertise to leverage LLMs effectively, and they frequently encounter writing tasks such as papers, assignments, and examinations. ChatGPT, in particular, has gained widespread popularity among college students worldwide. Consequently, universities urgently need robust detectors to address this issue, which has driven continuous advancements in the field of detection technology.

Research on detectors in this field can be broadly categorized into two directions. The first approach involves expanding the machine text corpus and enhancing the detector's performance using diverse training data. The second approach focuses on designing novel detector structures to improve overall performance. Both directions have yielded notable results, with detectors showcasing good performance on limited test sets in their respective research papers.

The versatility of LLMs, including the GPT family, enables students to exploit various prompts for academic cheating, thereby undermining detectors' effectiveness. However, achieving high performance solely on limited test sets falls short of adequately addressing real-world challenges. There is a pressing need to evaluate the robustness of models across a broader range of prompts and test sets, an aspect that has been largely overlooked in existing studies.

043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083

Our paper makes three contributions:

- **Highlighting the insufficient robustness of existing detectors through the example of academic writing cheating:** We demonstrate that solely adjusting the prompt is inadequate for ensuring the robustness of current detectors, particularly in the context of academic writing cheating.
- **Introducing a new detection approach for academic writing cheating using a Siamese network:** We analyze the academic writing cheating scenario and propose a novel detection approach based on a Siamese network. Our model exhibits superior prompt generalization capabilities compared to existing detectors, effectively addressing the issue of insufficient robustness when confronted with specific prompts.
- **Exploring the prompt-induced lack of robustness and evaluating model applicability:** We put forward a hypothesis to explain the reasons behind the lack of robustness caused by the prompt and provide evidence to support our claims. Furthermore, we demonstrate the broad applicability of our model based on this hypothesis.

The rest of the paper is organized as follows: Firstly, we review the literature and evaluate the existing detector’s lack of robustness in detecting academic cheating. Subsequently, we conduct an in-depth analysis of the academic cheating scenario, leading us to propose a new network specifically designed to address the robustness issue. Finally, we put forward a hypothesis regarding the factors contributing to the lack of robustness in generated articles, supported by our experimental evidence.

## 2 Related Work

With the popularity of LLM, many studies have explored the security problems that LLM may bring in recent years. Evan et al. conducted a comprehensive investigation into the potential security issues posed by LLMs and provided an overview of existing detection systems (Crothers et al., 2023). Stiff et al. analyzed the possible disinformation of false texts, tested the text on multiple platforms using the RoBERTa model, and analyzed whether the existing detection technology can detect the existing

generated text (Stiff and Johansson, 2022). Greshake, et al. pointed out that many applications now integrate LLMs as part of their functions (Greshake et al., 2023). However, LLMs may be affected by the input. If an attacker designs malicious input to mislead LLMs, it will likely cause data leakage and other security problems.

Researchers in detector development have explored strategies to optimize the training set for improved model performance. Notably, Liyanage et al. pioneered an AI-generated academic dataset using GPT-2, although it is considered inferior to the more advanced ChatGPT model currently available (Liyanage et al., 2022). Yuan et al. proposed BERTscore, a novel evaluation method for filtering high-quality generated text that closely resembles human writing (Yuan et al., 2021). Such text can be incorporated into the training set, thereby enhancing the performance of the detectors.

Researchers have also focused on optimizing the model itself. Jawahar et al. addressed the challenge of hybrid text, introducing a method to detect the boundary between machine-generated and human-written content, rather than solely distinguishing between the two (Jawahar et al., 2020). Zhao et al. conducted a comprehensive survey of various LLMs, analyzing their performance across multiple dimensions, including pre-training, adaptation tuning, utilization, and capacity evaluation. They also identified potential future development directions for LLMs (Zhao et al., 2023). Additionally, Mitchell et al. proposed a novel model utilizing a curvature-based criterion to determine whether a given passage was generated by an LLM (Mitchell et al., 2023).

Studies have also examined the robustness of detectors. Rodriguez et al. investigated the impact of dataset domain on detector performance, highlighting a significant decrease in performance when the training and test datasets differ in domain (Rodriguez et al., 2022). Their findings emphasized how the diversity of training sets directly affects the detector’s performance. Pu et al. analyzed the issue of insufficient robustness in existing detection systems by exploring changes in decoding or text sampling strategies (Pu et al., 2022). While previous research focused on robustness in terms of dataset domains and generative models’ parameters, this study highlights that prompt adjustments alone can significantly affect the robustness of the detector, particularly in the context of academic

**Simple prompt:** Write an abstract for a paper.

**Specific prompt:** Write an abstract for a paper about



Human-written paper's title
Universal Metrics for Large-scale Performance Analysis of Deep Neural Network.
Algorithms and Complexity of Range Clustering.
⋮

Figure 1: Examples of a simple prompt and a specific prompt.

cheating. The subsequent section will provide a demonstration of this phenomenon.

### 3 Asserting the Limitation of Existing Detectors

We conducted a simple preliminary test to highlight the prompt-induced limitations of a state-of-the-art AI-generated text detector.

#### 3.1 Dataset Construction

Since the release of GPT-3, OpenAI has allowed users to provide input prompts to shape the output text, enabling a wide range of functionalities. This inclusion of prompts significantly enhances the variation in generated text, presenting a more significant challenge for detection tasks. In contrast, the previous model, GPT-2, lacks prompt functionality and is irrelevant to the robustness of prompt-related issues. ChatGPT, a question-answering platform, does not offer APIs or adjustable parameters, making it unsuitable for generating large-scale datasets with diverse outputs. Hence, this paper uses GPT-3 for dataset generation, serving as the benchmark for our measurements. It is essential to clarify that throughout this paper, the term “GPT model” specifically refers to GPT-3.

For the human-written part of the dataset, we obtained the real human paper abstracts by collecting 500 samples from the arXiv dataset (Clement et al., 2019), which is available on Kaggle<sup>1</sup> and covers various fields. To create the AI-generated part of the dataset, we divided it into two subsets as depicted in Fig. 1. The “Simple prompt” subset consists of 500 GPT abstracts generated by GPT-3 using the prompt “Write an abstract for a professional paper.” The “Specific prompt” subset includes 500 GPT abstracts generated by GPT-3 using prompts beginning with "Write an abstract for a paper about" followed by the corresponding titles from the real human abstracts.

<sup>1</sup><https://www.kaggle.com/datasets/Cornell-University/arxiv>

#### 3.2 Detector Benchmark

Among the state-of-the-art detectors available, such as ChatGPT detector and GPTZero, many lack associated published articles or datasets for reproduction. Moreover, a significant number of these detectors do not provide APIs, making it impossible to conduct batch-testing experiments. Consequently, we have chosen the RoBERTa base OpenAI Detector (OpenAI detector in short) on Hugging Face<sup>2</sup>, a single-input binary classifier, as our target detector due to its availability and usability.

The detector demonstrates an impressive accuracy of 98% in detecting abstracts generated by **simple prompts** and 98% in identifying human-written abstracts. However, when it comes to abstracts generated by **specific prompts**, the accuracy rate **drops to only 87%**. This substantial reduction in performance by simply adding a human-written sentence to the prompt clearly indicates the limited robustness of existing detectors. Notably, specific prompts are commonly used in academic cheating scenarios, where students tailor their assignments or reports to meet specific requirements provided by their professors, utilizing prompts similar to the specific prompts used in this study. An example of the abstract generated using the corresponding title that was misclassified by the OpenAI detector is shown in Tab. 4 in the Appendix.

### 4 Our Solution

Our solution consists of two key components. Firstly, we analyze potential academic cheating scenarios and develop a **cheating model** specifically tailored to address these instances of cheating. Secondly, we propose a **novel detection system** designed to identify instances of academic cheating based on our developed model.

<sup>2</sup><https://huggingface.co/roberta-base-openai-detector>

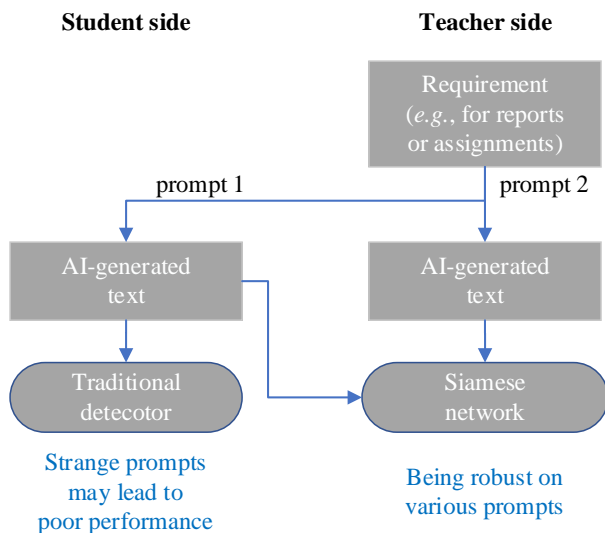


Figure 2: The proposed student cheating model.

#### 4.1 Student Cheating Model

As depicted in Fig.3, the model comprises two parties: the student side and the teacher side. Initially, the teacher assigns specific requirements for an academic task. Subsequently, a potentially deceitful student utilizes these requirements as input to generate an article using a generative model, such as GPT-3. The student may customize the provided requirements to evade detection, as discussed in Section3 with specific prompts. On the other hand, the teacher also proactively employs the generative model to generate an article. Then, the teacher uses a model to compare the similarities between the student’s submission and their own generated text in terms of content and style to determine whether the student engaged in cheating.

This cheating model closely resembles real-life situations where students’ assignments or examination articles are typically centered around specific topics and come with detailed requirements from teachers. To meet these requirements, students generally use the teacher’s instructions as input for generating their articles. Any slight modifications to the requirements or using different seeds for the generative model have minimal impact on the cheating model.

#### 4.2 Detection System

The network structure, as depicted in Fig. 3, involves the input of two articles:  $\mathbf{x}$  and  $\mathbf{y}$ . The article  $\mathbf{y}$  represents the teacher’s AI-generated article, while  $\mathbf{x}$  can either be a human-written article or an AI-generated one submitted by the student.

Our detector employs a pre-trained BERT net-

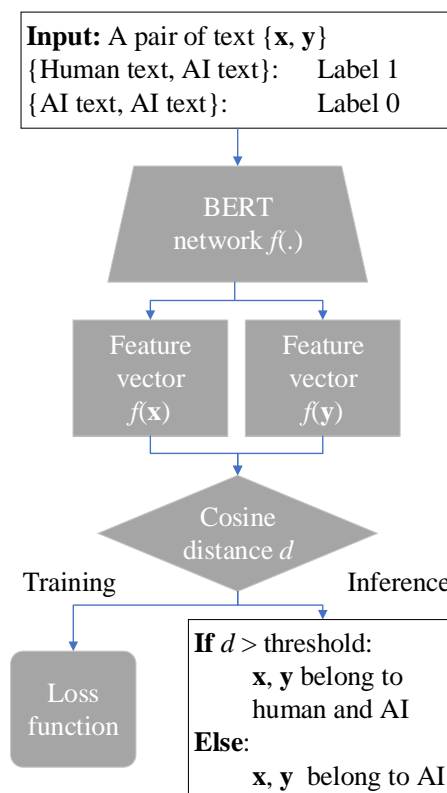


Figure 3: Overview of the proposed detector network.

work as a feature extractor, denoted as  $f(\cdot)$ , which is initialized with pre-trained weights. We fine-tune it using a supervised training approach. During the labeling of training data, if both  $\mathbf{x}$  and  $\mathbf{y}$  represent AI-generated articles, the label  $l$  is assigned as 0. Conversely, if  $\mathbf{x}$  corresponds to a human-written article and  $\mathbf{y}$  represents an AI-generated article, the label  $l$  is set as 1.

We use cosine distance  $\delta(\cdot, \cdot)$  for measuring the similarity between two feature vectors  $\mathbf{f}_x = f(\mathbf{x})$  and  $\mathbf{f}_y = f(\mathbf{y})$ , described in Eq. 1.

$$\delta(\mathbf{f}_x, \mathbf{f}_y) = 1 - \frac{\mathbf{f}_x \cdot \mathbf{f}_y}{\|\mathbf{f}_x\|_2 \|\mathbf{f}_y\|_2} \quad (1)$$

The loss function utilized during training is described by Eq. 2.

$$\mathcal{L} = l\delta(\mathbf{f}_x, \mathbf{f}_y)^2 + (1 - l)(2 - \delta(\mathbf{f}_x, \mathbf{f}_y))^2 \quad (2)$$

During the inference phase, our model calculates the cosine distance between the two input texts. A smaller distance indicates a higher similarity between  $\mathbf{x}$  and  $\mathbf{y}$ . As  $\mathbf{y}$  represents AI-generated text, a smaller distance suggests that  $\mathbf{x}$  is more likely to be generated by AI. Conversely,  $\mathbf{x}$  is more likely to be written by a human or contain a significant human contribution.

Table 1: Accuracy of the detectors on the prompt-generalization test set with level-n prompts.

Prompt variant	Prompt content	OpenAI detector (original)	OpenAI detector (fine-tuned)	Proposed detector
Human text		100%	98%	92%
Directly use requirement	Write an abstract for a paper about X	11%	35%	85%
Another expression	If you are a student, please complete the abstract of the article assigned by the teacher with topic X.	17%	46%	71%
Double GPT	Revise X then write an abstract about the revised text.	7%	15%	81%
Many $\rightarrow$ one	Find five human abstracts about X then summarize them into one.	11%	14%	81%

## 5 Experiment Results and Discussions

### 5.1 Experimental Design

We conducted experiments following similar settings as described in Section 3, but with an expanded dataset as outlined below:

- We utilized various levels of **specific prompts** in four different variants, as illustrated in Fig. 4 (and exemplified in Tab. 5 in the Appendix).
- The **training set** consisted of 2,000 human-written abstracts and 4,000 GPT-3 generated texts using **level-1** specific prompts. This dataset was employed for fine-tuning the detectors.
- For the **prompt-generalization test set**, we selected 100 human-written abstracts and generated 100 abstracts per each prompt variant that mimics different manipulative behaviors students may employ with **level n**. The prompt variants include “*Directly use requirement*,” which is the specific prompt we designed before. “*Another expression*” is where the student expresses the meaning of the requirement using different wording. The “*Double GPT*” variant involves using the generative model (GPT) twice, where the student modifies the original human idea X using GPT before generating the article. Lastly, the “*Many  $\rightarrow$  one*” variant simulates a common plagiarism method where the student collects five human articles about human idea X and combines them into a new article. These prompt

variants allow us to evaluate the detector’s performance in detecting different manipulative strategies employed by students. Examples of each variant are shown in Fig. 4 and Tab. 1.

- We extended the prompt-generalization test set to form the **human-contribution test set** by incorporating different levels of human contribution. Each variant in the test set represents a different level of human’s involvement in the generated text. The levels range from including only the field name, to including the title, summary of the abstract, and finally the entire abstract, denoted as 0, 1, 2, and n, respectively. By incorporating varying degrees of human contribution, we aim to assess the detector’s ability to distinguish between AI-generated text with different levels of human involvement. Examples of each level are shown in Fig. 4 (and Tab. 5 in the Appendix).
- For the **domain-generalization test set**, we chose 50 human-written abstracts and generated 50 abstracts per each generative model comprising OpenAI’s GPT-3, Perplexity’s customized GPT-3.5<sup>3</sup>, and the Falcon-7B<sup>4</sup>. All abstracts were generated using **level-1** and **level-2** prompts. This test set enables us to assess the detectors’ ability to generalize across different generative models, providing insights into their performance and adaptability in diverse AI-generated text scenarios.

<sup>3</sup><https://www.perplexity.ai/>

<sup>4</sup><https://falconllm.tii.ae/>

Table 2: Accuracy of the original OpenAI detector (before fine-tuning) in different levels. X denotes the human-written content incorporated into the prompts.

X level	Directly use requirement	Another expression	Double GPT	Many → one
level 0 (X = Field name)	100%	100%	99%	86%
level 1 (X = Title)	70%	74%	53%	72%
level 2 (X = Summary of abstract)	34%	24%	20%	29%
level n (X = Entire abstract)	11%	17%	7%	11%

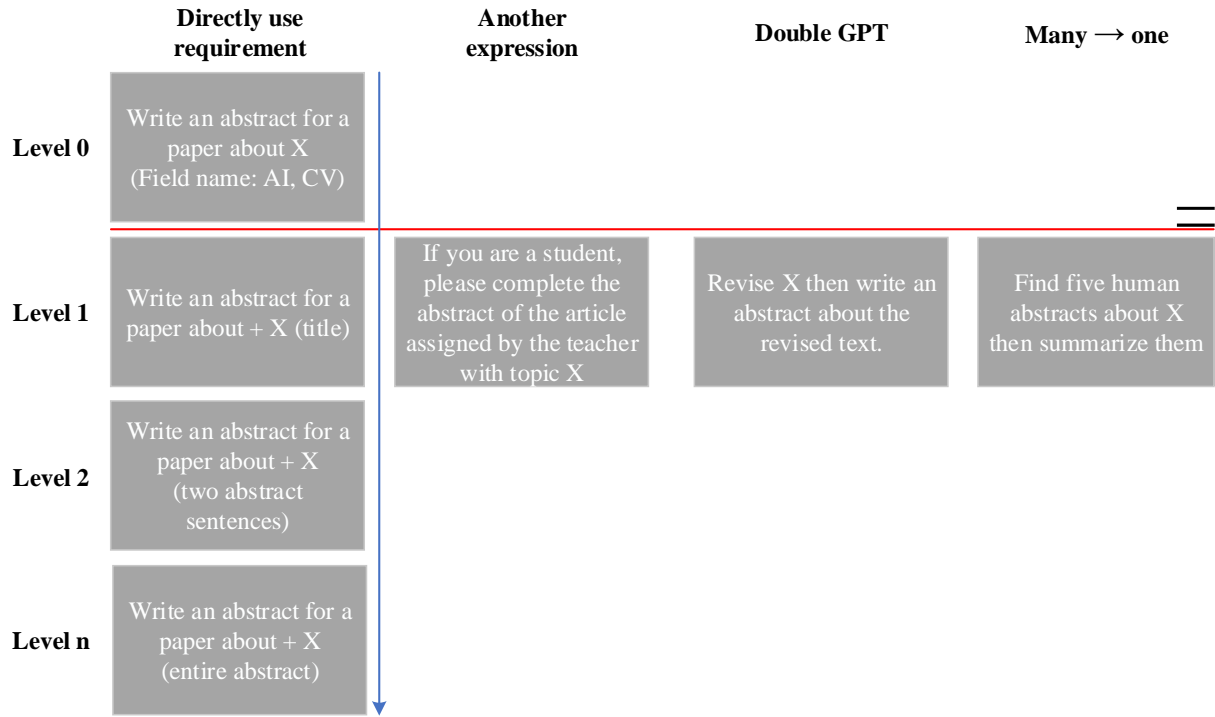


Figure 4: The prompts can be categorized into different levels based on the degree of human-written content. The horizontal red line indicates that prompts within the same level share similar characteristics. The vertical blue arrow illustrates that the generated articles become more challenging to classify accurately as the level increases.

374 Regarding the classification threshold (cosine  
 375 distance) employed by our detector, we have empirically  
 376 set it at 0.8. This threshold strikes a balance  
 377 between the false rejection rate and the false accep-  
 378 tance rate across various scenarios. However, it is  
 379 important to note that users have the flexibility to  
 380 adjust this threshold based on their individual use  
 381 cases and specific requirements.

## 382 5.2 Prompt-Variant Generalizability

383 We utilized the prompt-generalization test set to  
 384 assess the detectors’ performance in detecting vari-  
 385 ous prompt variants. As presented in Table 1, the  
 386 OpenAI detector exhibited a significant drop in per-  
 387 formance on different variants of prompts level n  
 388 (the most extreme cases), even after fine-tuning,  
 389 with a maximum true positive rate (TPR) of only

390 46%. In contrast, our model demonstrated superior  
 391 generalizability, achieving a minimum TPR of 71%  
 392 on the “another expression” specific prompts. This  
 393 implies that in academic cheating scenarios, our  
 394 model can effectively detect the usage of GPT by  
 395 students, regardless of the complexity of the pro-  
 396 fessor’s requirements and the inclusion of a certain  
 397 amount of human-written content in the prompts  
 398 (approximately 200 to 250 words as an abstract).

399 In terms of the true negative rate (TNR), which  
 400 evaluates the detectors’ capability to accurately  
 401 identify human-written text, our detector achieved  
 402 a commendable accuracy of 92%. Although this  
 403 is slightly lower than the fine-tuned OpenAI de-  
 404 tector (98%) and its original version (100%), it is  
 405 a reasonable trade-off considering the decrease in  
 406 the TPRs of the OpenAI detector. Furthermore,

Table 3: Accuracy (or TPR) of the detectors on the text generated by different LLMs. OpenAI detector, a binary classifier, only needs one input. Our detector, besides the query text, requires the corresponding generated text (from the teacher) as an anchor. Within each cell, the upper number represents the result on level-1 prompts, while the lower number represents the result on level-2 prompts.

Source of input text	OpenAI detector (original)	OpenAI detector (fine-tuned)	Proposed detector		
			GPT-3 text as anchor	Falcon-7B text as anchor	Perplexity text as anchor
Human	100%	98%	92%	70%	90%
	100%	98%	92%	12%	90%
GPT-3	70%	99%	95%		
	35%	98%	100%		
Falcon-7B	16%	92%	60%	70%	
	13%	57%	12%	96%	
Perplexity	47%	98%	100%		100%
	53%	98%	70%		100%

users have the flexibility to adjust the classification threshold according to their specific use cases and requirements.

To investigate the drop in the OpenAI detectors’ performance, we examined the impact of reducing human contribution in prompts using the human-contribution test set. Results in Table 2 showed that the original OpenAI detector ideally detected AI-generated text with level-0 prompts, except for “many → one” prompts (86% accuracy). Performance remained acceptable at level 1 but deteriorated significantly at level 2 and beyond. This is unacceptable in real-life scenarios where malicious students may strategically add additional keywords or phrases to make their generated text more convincing and harder to detect.

### 5.3 Domain Generalizability

Although GPT has become mainstream, students may utilize several other text-generation models based on LLMs to avoid detection. To assess the detectors’ effectiveness, we conducted tests using the domain-generalization test set. It is important to note that all detectors were fine-tuned solely using GPT-3 generated text.

The results are presented in Table 3. The original OpenAI detectors struggled to perform effectively in most cases, while its fine-tuned version achieved the highest accuracies except when dealing with text generated by Falcon-7B using level-2 prompts. Our proposed detector performed highly on the GPT variants (GPT-3 and customized GPT-3.5). However, it showed limited generalizability when faced with Falcon’s generated text using anchor text generated by other LLMs.

We hypothesized that during the training of our Siamese-based detector with the proposed cheating model, the detector learned to identify authorship information. It distinguished GPT as one author and humans as another. When a new “author” (Falcon-7B) emerged, the detector struggled to assign its text to either the human or GPT. When using asymmetric pairs as input, the scores fell around the decision threshold, leading to degraded performance. Conversely, when using pairs of Falcon-7B’s text, the detector treated them as originating from the same author, resulting in improved accuracy.

To improve the inter-model generalizability of our detector, teachers can select representative models from popular LLM families such as GPT, LLaMA (Touvron et al., 2023), and Falcon to create multiple anchor texts for multiple comparisons.

## 6 Hypothesis for the Prompt-Induced Lack of Robustness

As demonstrated in the previous section, traditional detectors exhibit limited robustness due to the infinite possibilities of prompts. While we evaluated specific prompts related to academic cheating, it is crucial to acknowledge that the prompts we examined cannot encompass the entire spectrum of academic cheating scenarios. To systematically address this issue, we generalized the result in Tab.2 to form a hypothesis that aims to (1) illuminate the potential factors underlying the reduced robustness of traditional detectors and (2) substantiate the generalizability of our chosen of prompts.

Our hypothesis can be illustrated using Figure 4

474	and can be explained as follows:		520
475			521
476			522
477			523
478			524
479			525
480			526
481			527
482			528
483			529
484			530
485			531
486			532
487			533
488			534
489			535
490			536
491			537
492			538
493			539
494			540
495			541
496			542
497			543
498			544
499			545
500			546
501			547
502			548
503			549
504			550
505			551
506			552
507			553
508			554
509			555
510			556
511			557
512			558
513			559
514			560
515			561
516			562
517			563
518			564
519			565
			566
			567



568	Colin B Clement, Matthew Bierbaum, Kevin P O’Keeffe, and Alexander A Alemi. 2019. On the use of arxiv as a dataset. <i>arXiv preprint arXiv:1905.00075</i> .	Juan Rodriguez, Todd Hay, David Gros, Zain Shamsi, and Ravi Srinivasan. 2022. <b>Cross-domain detection of GPT-2-generated technical text</b> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1213–1233, Seattle, United States. Association for Computational Linguistics.	620 621 622 623 624 625 626 627
572	Evan Crothers, Nathalie Japkowicz, and Herna Viktor. 2023. <b>Machine generated text: A comprehensive survey of threat models and detection methods</b> .		
575	Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. 2023. Not what you’ve signed up for: Compromising real-world llm-integrated applications with indirect prompt injection.	Harald Stiff and Fredrik Johansson. 2022. <b>Detecting computer-generated disinformation</b> . <i>International Journal of Data Science and Analytics</i> , 13.	628 629 630
576			
577			
578			
579			
580	Ganesh Jawahar, Muhammad Abdul-Mageed, and Laks V. S. Lakshmanan. 2020. <b>Automatic detection of machine generated text: A critical survey</b> .	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	631 632 633 634 635 636
581			
582			
583	Keenan Jones, Enes Altuncu, Virginia N. L. Franqueira, Yichao Wang, and Shujun Li. 2022. <b>A comprehensive survey of natural language generation advances from the perspective of digital deception</b> .	Haoran Wu, Wenxuan Wang, Yuxuan Wan, Wenxiang Jiao, and Michael Lyu. 2023. <b>Chatgpt or grammarly? evaluating chatgpt on grammatical error correction benchmark</b> .	637 638 639 640
584			
585			
586			
587	M. Kestemont, E. Manjavacas, I. Markov, J. Bevendorff, M. Wiegmann, E. Stamatatos, B. Stein, and M. Potthast. 2021. Overview of the cross-domain authorship verification task at pan 2021. In <i>CLEF-WN 2021 - Proceedings of the Working Notes of CLEF 2021 - Conference and Labs of the Evaluation Forum</i> , volume 2936 of <i>CEUR Workshop Proceedings</i> , pages 1743–1759. CEUR-WS.	Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. <b>BARTScore: Evaluating generated text as text generation</b> .	641 642 643
588			
589			
590			
591			
592			
593			
594			
595	Cyril Labbé and Dominique Labbé. 2012. <b>Duplicate and fake publications in the scientific literature: how many SCiGen papers in computer science?</b> <i>Scientometrics</i> , pages 10.1007/s11192–012–0781–y.	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. <b>A survey of large language models</b> .	644 645 646 647 648 649 650
596			
597			
598			
599	Vijini Liyanage, Davide Buscaldi, and Adeline Nazarenko. 2022. <b>A benchmark corpus for the detection of automatically generated text in academic publications</b> .	<b>A Examples of Human’s and AI’s Texts and Their Detection.</b>	651 652
600			
601			
602			
603	Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. 2023. <b>Detectgpt: Zero-shot machine-generated text detection using probability curvature</b> .	Tab. 4 shows an example of the abstract generated using the corresponding title that was misclassified by the OpenAI detector. Tab 5 shows examples of prompts with different levels of human-written contents (factor X).	653 654 655 656 657
604			
605			
606			
607	OpenAI. 2023. GPT-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .	<b>B Scientific Artifacts Detail</b>	658
608			
609	Jiameng Pu, Zain Sarwar, Sifat Muhammad Abdullah, Abdullah Rehman, Yoonjin Kim, Parantapa Bhat-tacharya, Mobin Javed, and Bimal Viswanath. 2022. <b>Deepfake text detection: Limitations and opportunities</b> .	This paper adheres strictly to OPENAI’s terms of use <sup>5</sup> , and no violations have occurred. It is important to note that OPENAI has not provided specific guidelines regarding expected model application scenarios. The human texts used in this paper were sourced from the publicly available arXiv dataset (Clement et al., 2019). No personally identifiable information has been included in the study.	659 660 661 662 663 664 665 666
610			
611			
612			
613			
614	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. " why should i trust you?" explaining the predictions of any classifier. In <i>Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining</i> , pages 1135–1144.		
615			
616			
617			
618			
619			

<sup>5</sup><https://openai.com/policies/terms-of-use>

Table 4: An example that OpenAI detector misclassified a GPT-generated text as written by humans. Leveraging LIME (Ribeiro et al., 2016), we demonstrate that the OpenAI Detector exhibited high confidence in its classification of the [blue sentences](#) as human-written. In contrast, our detector accurately distinguishes between human-generated text and GPT-generated text, correctly classifying both with precision.

GPT Text	Human Text
<p>This paper presents a new knowledge selection method for knowledge-grounded conversation generation. <a href="#">This method, called Difference-aware Knowledge Selection (DKS)</a>, leverages the difference between a given conversation context and the associated knowledge to determine the most relevant knowledge to use. DKS first computes the semantic similarity between the conversation context and the available knowledge. It then uses a reinforcement learning algorithm to select the knowledge with the highest reward, which is calculated by the semantic similarity and the expected conversation turn difference. <a href="#">Experimental results demonstrate that the DKS method outperforms baseline methods in terms of both response quality and diversity.</a></p>	<p>In a multi-turn knowledge-grounded dialog, the difference between the knowledge selected at different turns usually provides potential clues to knowledge selection, which has been largely neglected in previous research. In this paper, we propose a difference-aware knowledge selection method. It first computes the difference between the candidate knowledge sentences provided at the current turn and those chosen in the previous turns. Then, the differential information is fused with or disentangled from the contextual information to facilitate final knowledge selection. Automatic, human observational, and interactive evaluation shows that our method is able to select knowledge more accurately and generate more informative responses, significantly outperforming the state-of-the-art baselines. The codes are available at <a href="https://github.com/chujiezheng/DiffKS">https://github.com/chujiezheng/DiffKS</a>.</p>

Table 5: Examples of human-written contents from different levels of prompts. It is obvious that as the level increases, the length and complexity of X increase. Therefore, it can be considered that the higher the level, the more human ideas X contains.

Level 0	Level 1	Level 2	Level n
Field name	Title	Summary of abstracts	Whole abstract
AI	DeepStruct: Pretraining of Language Models for Structure Prediction	We introduce a method for improving the structural understanding abilities of language models. Unlike previous approaches that finetune the models with task-specific augmentation, we pretrain language models on a collection of task-agnostic corpora to generate structures from text. Our structure pretraining enables zero-shot transfer of the learned knowledge that models have about the structure tasks.	We introduce a method for improving the structural understanding abilities of language models. Unlike previous approaches that finetune the models with task-specific augmentation, we pretrain language models on a collection of task-agnostic corpora to generate structures from text. Our structure pretraining enables zero-shot transfer of the learned knowledge that models have about the structure tasks. We study the performance of this approach on 28 datasets, spanning 10 structure prediction tasks including open information extraction, joint entity and relation extraction, named entity recognition, relation classification, semantic role labeling, event extraction, coreference resolution, factual probe, intent detection, and dialogue state tracking. We further enhance the pretraining with the task-specific training sets. We show that a 10B parameter language model transfers non-trivially to most tasks and obtains state-of-the-art performance on 21 of 28 datasets that we evaluate.

667

## **C Computational Experiment Detail**

668

669

670

671

672

Our proposed model contained 108.57M parameters and was trained for two hours on a single NVIDIA A100 GPU. To ensure the reliability of the results, we conducted two runs and averaged the outcomes reported in this paper.