

# MGTBench: Benchmarking Machine-Generated Text Detection

Xinlei He<sup>†</sup> Xinyue Shen<sup>†</sup> Zeyuan Chen<sup>‡</sup> Michael Backes<sup>†</sup> Yang Zhang<sup>†</sup>

<sup>†</sup>*CISPA Helmholtz Center for Information Security* <sup>‡</sup>*Individual Researcher*

## Abstract

Nowadays, powerful large language models (LLMs) such as ChatGPT have demonstrated revolutionary power in a variety of natural language processing (NLP) tasks such as text classification, sentiment analysis, language translation, and question-answering. Consequently, the detection of machine-generated texts (MGTs) is becoming increasingly crucial as LLMs become more advanced and prevalent. These models have the ability to generate human-like language, making it challenging to discern whether a text is authored by a human or a machine. This raises concerns regarding authenticity, accountability, and potential bias. However, existing methods for detecting MGTs are evaluated using different model architectures, datasets, and experimental settings, resulting in a lack of a comprehensive evaluation framework that encompasses various methodologies. Furthermore, it remains unclear how existing detection methods would perform against powerful LLMs.

In this paper, we fill this gap by proposing the first benchmark framework for MGT detection against powerful LLMs, named MGTBench. Extensive evaluations on public datasets with curated answers generated by various powerful LLMs such as ChatGPT and GPT4All show that most of the current detection methods perform less satisfactorily against MGTs except ChatGPT Detector and LM Detector. Our ablation study shows that model-based detection methods can achieve similar performance with much fewer training samples and can transfer well to other datasets/LLMs. Furthermore, we delve into a more challenging task: text attribution, where the goal is to identify the originating model of a given text, i.e., whether it is a specific LLM or authored by a human. Our findings indicate that the LM Detector performs the best in the text attribution task. However, it is worth noting that only a small fraction of adversarial-crafted perturbations on MGTs might evade the LM Detector, underscoring the need for more robust MGT detection methods. We envision that MGTBench will serve as a benchmark tool to accelerate future investigations involving the evaluation of powerful MGT detection methods on their respective datasets and the development of more advanced MGT detection methods.<sup>1</sup>

## 1 Introduction

Large language models (LLMs), such as T5 [25], GPT3 [10], PaLM [11], and more powerful LLMs such as ChatGPT [6], have been a significant breakthrough in the field of natural language processing (NLP). With huge numbers of parameters and being trained with massive amounts of data, LLMs have shown remarkable performance in various real-world applications, e.g., education, customer service, and finance.

While powerful LLMs have shown impressive qualities in terms of achieving remarkable performance in various tasks and generating human-like texts, several limitations and ethical concerns have to be taken into account. First, LLMs may generate text that sounds realistic but may not be entirely accurate or factual [19]. Second, the misuse of LLMs may raise concerns in education, making a fair judgment impossible [30]. Also, it is difficult to trace back the machine-generated text to its source, which raises concerns about accountability, especially when the content is used to spread misinformation or propaganda [34].

To address these issues, researchers have considered automatic detection methods that can identify the machine-generated text (MGT) from the human-written text (HWT). Concretely, those methods can be concluded into two categories, i.e., metric-based methods and model-based methods. For the metric-based methods [15, 21, 28], metrics such as log-likelihood, word rank, and predicted distribution entropy are used to determine whether a text belongs to MGT or HWT. Regarding the model-based methods [4, 16, 28, 35], classification models are trained by both MGTs and HWTs.

Overall, existing MGT detection methods have been studied under different LLMs and different datasets, albeit in isolation. Also, powerful LLMs such as ChatGPT demonstrate remarkable performance across diverse tasks, providing high-quality, human-like answers. This prompts the need for a holistic evaluation of these methods against powerful LLMs. To do this, we develop a comprehensive benchmark of MGT detection, namely MGTBench. MGTBench follows a modular design, consisting of the input module, the detection module, and the evaluation module. Until now, we have implemented ten MGT detection methods. We also take advantage of the powerful LLMs including ChatGPT, ChatGLM, Dolly, ChatGPT-turbo, GPT4All, and StableLM, which are the most popular and powerful LLMs (thus far), to produce MGTs on existing HWT datasets. MGTBench enables researchers to benchmark various methods on differ-

<sup>1</sup>Source code and datasets are available at <https://github.com/xinleihe/MGTBench>.

ent datasets. Its modular design facilitates the integration of additional detection methods, as well as the plugging in of datasets and models.

**Evaluation.** We perform an extensive measurement study over the ten detection methods against six LLMs. Our measurement is performed on standard benchmark question-answering (QA) datasets including TruthfulQA [19], SQuAD1 [26], and NarrativeQA [18]. Note that for each dataset, besides the answers generated by humans (labeled as HWTs), we also query each LLM with the questions (and the contexts if they are included on the dataset) to obtain the answers from the LLM (labeled as MGTs). In this case, we label the three datasets as TruthfulQA, SQuAD1, and NarrativeQA, respectively (see Section 4 for more details).

Extensive evaluations show that the LM Detector outperforms other detection methods and reaches the best performance. For instance, to differentiate whether the texts are generated by humans or ChatGPT on TruthfulQA, the detection F1-score is 0.997 with LM Detector while only 0.945 with Log-Likelihood. Also, our ablation study shows that the answers generated by LLMs are much longer than human-written answers. To reduce the potential bias caused by the number of words in a text, we perform extra experiments with texts within 25 words and observe a clear degradation in the detection performance. For instance, after filtering the texts, the detection AUC of Log-Rank drops from 0.905 to 0.670 on TruthfulQA when the task is differentiating texts generated by ChatGPT-turbo or humans. On the other hand, ChatGPT Detector and LM Detector consistently exhibit high performance. Also, we observe that the model-based methods can easily adapt to new datasets with only limited training samples. For instance, with 10 training samples, the detection F1-score reaches 0.974 and 0.990 for ChatGPT Detector with ChatGPT and ChatGLM as the LLMs. Additionally, we find that the LM Detector can transfer well across different datasets/LLMs. Later, we consider a more challenging task, i.e., text attribution, where the goal is to identify the exact model to generate the text, e.g., either human or one of the LLMs. Our evaluation shows that the model-based detection methods perform better than the metric-based detection method in the text attribution task. For example, on TruthfulQA, the LM Detector achieves 0.631 F1-score while the F1-score is only 0.134 and 0.255 for Rank and GLTR.

We then investigate whether the best-performing detection methods (ChatGPT Detector and LM Detector) are robust against adversarial attacks. We empirically show that LM Detector is more stable than ChatGPT Detector, but still vulnerable to the adversarial perturbation introduced to the texts (see Table 5 for more details).

In summary, we make the following contributions:

- We propose MGTBench, a benchmarking framework for MGT detection against powerful LLMs.
- Our empirical evaluation shows that the LM Detector outperforms other detection methods in both MGT detection and text attribution tasks. Furthermore, we find that the LM Detector exhibits strong generalization capabilities when applied to different datasets and LLMs.

- While the LM Detector demonstrates greater robustness compared to the ChatGPT Detector, we also discover its vulnerability to adversarial attacks. This highlights the need for future research to develop more resilient MGT detection methods.

## 2 Preliminary and Related Work

### 2.1 Text Generation with Large Language Models

The recent advancements in LLMs can be traced back to the transformer architecture proposed by Vaswani et al. [33], which introduces the self-attention mechanisms to allow the model to focus on different regions of the input sequence. Later Radford et al. [23] develop the generative pre-trained transformer (GPT) based on the transformer architectures. Being trained with a large corpus of text data, GPT achieves superior performance on a wide range of language generation tasks. GPT2 [24], a larger version of GPT that contains more parameters and is trained on a larger corpus, is developed and achieves better performance than GPT. GPT3 [10] is the third generation of GPT with over 175B parameters, which has been shown to be capable of generating coherent and contextually appropriate text even in situations where it is provided with minimal input or guidance. Since November 2022, OpenAI releases ChatGPT [6], which is trained based on the GPT-3.5 architectures and leverages Reinforcement Learning from Human Feedback (RLHF) [12, 29] to improve its generation ability. ChatGPT shows revolutionary capabilities in generating coherent and relevant texts, which can be integrated into various applications, such as chatbots, customer service, and education.

Another notable model series in the field is the masked language model. Bidirectional Encoder Representations from Transformers (BERT) developed by Devlin et al. [13] is one of the most representative models that is pre-trained using the masked language modeling task. Liu et al. [20] develop RoBERTa, which leverages a robustly optimized BERT pre-training approach and can reach an even better performance than BERT in various tasks. With the widespread use of LLMs for generating texts, concerns about authenticity, accountability, and potential bias have also been raised.

In this paper, we consider six representative powerful LLMs, including ChatGPT, ChatGPT-turbo, ChatGLM, Dolly, GPT4All, and StableLM.

**ChatGPT [6].** ChatGPT is an advanced large language model built upon the GPT-3.5 architecture and specifically designed to generate highly human-like responses. To achieve this capability, ChatGPT is fine-tuned with Reinforcement Learning from Human Feedback (RLHF) [29] where human trainers actively engage in conversations with ChatGPT. We refer to the model in March 1st, 2023 version as ChatGPT.

**ChatGPT-turbo [6].** In addition to the aforementioned model, we also consider the latest iteration of ChatGPT, namely ChatGPT-turbo to further investigate the impact of model versions on MGT Detectors.

**ChatGLM [2].** ChatGLM is another large language model based on the GLM (General Language Model) frameworks [14]. The training process involved a combination of supervised fine-tuning, feedback bootstrap, and RLHF, thus enabling it to generate responses that are in accordance with human-like patterns and preferences.

**Dolly [3].** Dolly is an instruction-following large language model released by Databricks. It is trained on around 15K instruction/response records, generated by Databricks employees in capability domains such as brainstorming, classification, closed QA, generation, information extraction, open QA, and summarization.

**GPT4All [7].** GPT4All is an open-source assistant-style large language model. It is trained over a massive curated corpus including word problems, story descriptions, multi-turn dialogue, and code.

**StableLM [5].** StableLM is an auto-regressive language model based on the NeoX transformer architecture [9]. It is fine-tuned on various chat and instruction-following datasets.

## 2.2 Machine-Generated Text Detection

To address the above-mentioned issues, researchers have developed various MGT detection methods [8, 15–17, 21, 22, 28, 31, 32, 35]. Current detection methods can be divided into two categories, i.e., metric-based methods and model-based methods. Generally speaking, metric-based methods leverage pre-trained LLMs to process the text and extract distinguishable features from it, e.g., the rank or entropy of each word in a text conditioned on the previous context. In this paper, we consider six metric-based detection methods, including Log-Likelihood, Rank, Log-Rank, Entropy, GLTR, and DetectGPT.

**Log-Likelihood [28].** This approach leverages a language model to measure the token-wise log probability. Concretely, given a text, we average the token-wise log probability of each word to generate a score for this text. Note that a larger score denotes the text is more likely to be machine-generated.

**Rank [15].** For each word in a text, given its previous context, we can calculate the absolute rank of this word. Then, for a given text, we compute the score of the text by averaging the rank value of each word. Note that a smaller score denotes the text is more likely to be machine-generated.

**Log-Rank [21].** Slightly different from the Rank metric that uses the absolute rank, the Log-Rank score is calculated by first applying the log function to the rank value of each word.

**Entropy [15].** Similar to the Rank score, the Entropy score of a text is calculated by averaging the entropy value of each word conditioned with its previous context. As mentioned by previous work [15, 21], the machine-generated text is more likely to have a lower Entropy score.

**GLTR [15].** GLTR is developed as a support tool to facilitate the labeling process of whether a text is machine-generated. In our evaluation, we follow the suggestion of Guo et al. [16] and consider the Test-2 features (i.e., the fraction of words that rank within 10, 100, 1,000, and others). Note that one can easily implement other sets of features in MGTBench.

**DetectGPT [21].** Mitchell et al. [21] propose DetectGPT that measures the change of the model’s log probability function by adding minor perturbation to the original text. The intuition is that the text derived from an LLM has a tendency to be in the local optimal of the model’s log probability function. Therefore, any minor perturbation of model-generated text tends to have a lower log probability under the model than the original text, while minor perturbation of human-written text may have a higher or lower log probability than the original text.

Regarding the model-based methods, a classification model is usually trained using a corpus that contains both HWTs and MGTs. By doing this, the classification model is expected to have the capability in identifying MGTs from a given corpus. In this paper, we consider four model-based methods.

**OpenAI Detector [28].** The OpenAI Detector is used to detect GPT2-generated output, which was created by fine-tuning a RoBERTa model using outputs from the largest GPT2 model (i.e., with 1.5B parameters). This model is capable of predicting whether a given text was machine-generated or not.

**ChatGPT Detector [16].** ChatGPT Detector is developed by Guo et al. [16] to distinguish human-written texts from ChatGPT-generated texts. The model was created by fine-tuning a RoBERTa model using the HC3 [16] dataset. The detector is based on the RoBERTa model. The authors provide two ways to train the RoBERTa model. The first one only leverages the pure answered text, and the second one leverages the question-answer text pair to jointly train the model. In our evaluation, we consider the first one to be consistent with other detection methods.

**GPTZero [4].** GPTZero is an MGT analyzer tool that uses two main measures, i.e., perplexity and burstiness, to determine whether a text is machine-generated or written by a human.<sup>2</sup> GPTZero provides the publicly accessible API to produce a confidence score about how likely a text is generated by the machine. Note that in our initial experiment, we find that directly leveraging 0.5 as the threshold might have a less satisfying performance. Therefore, given the training dataset that contains both HWTs and MGTs, we first calculate the average confident scores of HWTs and MGTs, respectively. Then, we average the two scores as the threshold to get a concrete prediction for each confident score. We name this method GPTZero-align.

**LM Detector.** Besides the previous methods, the detector can also be built by fine-tuning the pre-trained language model (LM) with an extra classification layer. Here we take the BERT model as an example to evaluate its efficacy following Ippolito et al. [17].

Compared to previous work [22, 32], our study offers notable advancements in the field of MGT detection by integrating more detection methods into our benchmarking framework. Moreover, we extend the evaluation to more powerful LLMs such as ChatGPT. This broader scope enables us to

<sup>2</sup><https://www.makeuseof.com/gptzero-detect-ai-generated-text/>.

**Table 1: The prompts we used to obtain answers from LLMs. Note that <context> and <question> are provided by the original dataset. The prompts are adopted from [1].**

| Dataset     | Prompt  |
|-------------|---|
| TruthfulQA  | <question>  |
| SQuAD1      | I will provide a passage and a question to you. The answer should be extracted from the context. You need to return me your answer. The passage is <context> and the question is <question>. Now, please answer the question. |
| NarrativeQA | I will provide a context and a question to you. You need to answer me the question based on the context. The context is: <context> The question is: <question>  |

gain deeper insights into the performance and robustness of various detection methods when applied to a wider range of LLMs. By incorporating these enhancements, our work significantly contributes to the existing literature and paves the way for further advancements in MGT detection research.

### 3 MGTBench

In this section, we introduce MGTBench, a modular framework designed to benchmark MGT detection methods. Currently, we have provided reference implementations of the six metric-based detection methods and easy-to-use APIs for the four model-based methods we mentioned before.

#### 3.1 Modular Design

MGTBench consists of three different modules, including the *input module*, *detection module*, and *evaluation module*.

**Input Module.** In the input module, we provide specific dataset pre-processing functions for different datasets and our code base is easy to cope with datasets from HuggingFace, which facilitates the future developments of different users.

**Detection Module.** This module implements different metric-based and model-based detection methods with a standardized input/output format. Currently, we support ten different detection methods.

**Evaluation Module.** This module is used to evaluate the performance of different detection methods. Now, we provide five different evaluation metrics, including accuracy, precision, recall, F1-score, and AUC, which are the commonly used metrics to evaluate classification performance.

#### 3.2 Using MGTBench

To be best of our knowledge, MGTBench is the most comprehensive benchmark tool for MGT detection against powerful LLMs. Users can leverage MGTBench on their own dataset for a comprehensive risk assessment of potential MGTs in the dataset. On the other hand, researchers can leverage MGTBench as a tool to evaluate new MGT detection/generation methods. As MGTBench follows a modular design, its input and evaluation modules can be easily re-used by new detection methods. Also, new detection methods can be easily implemented within the standardized API provided by MGTBench. Moreover, MGTBench integrates well with

HuggingFace, given the fact that many model-based detection methods have published or are willing to publish their models into HuggingFace, MGTBench can be seamlessly updated to fit the new model-based detection methods. MGTBench is under continuous development and we will include more detection methods as well as analysis tools in the future.

## 4 Experimental Settings

### 4.1 Datasets

**TruthfulQA.** This dataset is derived from the TruthfulQA [19] dataset that contains 817 questions from 38 categories including health, law, finance, and politics.

**SQuAD1.** This dataset is derived from the SQuAD1 [26] dataset that contains more than 100,000 question-answer pairs selected from more than 500 articles. We randomly sample 1,000 records from the SQuAD1 dataset as our evaluation dataset.

**NarrativeQA.** This dataset is derived from the NarrativeQA [18] dataset that contains stories and corresponding questions designed to test reading comprehension. We randomly sample 1,000 records from the NarrativeQA dataset as our evaluation dataset.

For each question, we additionally query different LLMs with the question (and the context if it is included in the dataset) and obtain the answers as MGTs for different LLMs. To avoid being affected by the chat history, we create a new conversation with each LLM for each question.

We show the prompt we used to query LLMs in Table 1. Note that for each entry in a dataset, we have the human answer and six LLM-generated answers (see Section 2.1 for more details about the six LLMs). We only keep entries with more than 1 word for human and LLM-generated answers. Then, we randomly split 80% of the entries as the training set and the rest as the testing set.

### 4.2 Tasks

In our evaluation, our primary focus is on the MGT detection task, which involves detecting whether a given text is generated by a human or a machine. To accomplish this, we establish a binary classification task for texts generated by humans and each LLM, such as Human vs. ChatGPT, Human vs. ChatGLM, and so on. Note that we also consider a



more complex task, namely text attribution (see Section 5.2). This task seeks to pinpoint the precise model responsible for generating the text. Essentially, when presented with a text, our goal is to determine whether it is generated by the human or one of the six LLMs. This task can be viewed as a seven-category classification scenario.

### 4.3 Detection Methods

For metric-based methods, we use GPT2-medium as the base model in our experiments since it can already reach good performance with limited cost. Given the metrics extracted with the GPT2-medium, we additionally build a logistic regression model on top of it to provide concrete predictions. For model-based methods, we directly use the pre-trained models, which can be obtained from HuggingFace.<sup>3</sup>. Concretely, for OpenAI Detector, we use the RoBERTa-base version of it as it usually gives better detection performance. For ChatGPT Detector, we leverage the provided RoBERTa-base model of it. For LM Detector, we leverage the distilled BERT-base model as it has superior performance and modest expenditure. Note that for OpenAI Detector and ChatGPT Detector, the pre-trained models are already optimized for the MGT detection task (binary classification), therefore, we do not further fine-tune them unless otherwise mentioned. It is also worth mentioning that, in the text attribution task, we fine-tune all detection methods as the number of classes increases from 2 to 7.

### 4.4 Evaluation Metrics

MGTBench supports various metrics to evaluate performance, including accuracy, precision, recall, F1-score, and AUC (area under the ROC curve). In our evaluation, unless otherwise mentions, we use F1-score as the main evaluation metric.

## 5 Evaluation

We first present the experiment results on the MGT detection task. As shown in Table 2, we observe that LM Detector achieves the best performance across different datasets/LLMs. For instance, on TruthfulQA, the LM detector achieves 0.997 F1-score against ChatGLM. This is expected as the LM detector is trained on the corpus with both HWTs and MGTs, and the LM’s ability in capturing the context and coherence information greatly aids in differentiating the pattern variations between these two types of texts. We also observe that the ChatGPT detector shows good performance in differentiating HWTs from MGTs, especially for the models from the GPT family, i.e., ChatGPT, ChatGPT-turbo, and GPT4All. For instance, on TruthfulQA, the F1-score is 0.974, 0.967, and 0.946 for ChatGPT, ChatGPT-turbo, and GPT4All, respectively. We suspect the reason is that the ChatGPT detector is trained on the corpus with both human answers and ChatGPT-generated answers, which can better capture the “machine” pattern on the GPT model family.

<sup>3</sup><https://huggingface.co/>.

Also, we observe that metric-based methods such as Log-Likelihood, Log-Rank, and Entropy have relatively good performances on different datasets as well. For instance, Log-Likelihood reaches 0.921, 0.736, and 0.725 F1-score on TruthfulQA, SQuAD1, and NarrativeQA when distinguishing ChatGPT-generated answers from HWTs. This implies that the MGT better fits the detection model’s “expectation” about what it suppose to be even if the detection model (GPT2-medium) is different from the one that is used to produce MGT (one of the six LLMs).

On the other hand, we find that the OpenAI Detector’s performance is not satisfying (e.g., 0.639 F1-score on TruthfulQA with ChatGPT-generated texts as MGTs). This might be credited to the fact that this detector is trained on MGTs produced by GPT2 while the MGTs produced by larger LLMs such as ChatGPT or Dolly have higher quality, which makes it harder to be detected.

Last, we also observe that GPTZero-align has less satisfying performance on detecting MGTs (e.g., 0.707 F1-score on NarrativeQA with ChatGPT). One possible reason is that GPTZero is more effective in detecting longer MGTs, which is also suggested on their official website [4].

**Detection Efficiency.** We then quantify the time cost of each detection method on different datasets. Here we consider ChatGPT as the LLM for a case study since other LLMs show similar time costs. As shown in Table 3, we can observe that most of the detection methods have similar time costs except the DetectGPT and GPTZero, which cost significantly higher time than the others. This is because DetectGPT requires multiple rounds of perturbation to the text to get a good estimation of the log probabilities’ change. And GPTZero needs to query the public API, where the internet latency has to be taken into account.

In general, we consider LM Detector and ChatGPT Detector as better detection methods as they achieve the best detection performance. Moreover, their associated time cost is relatively low, making them comparable to other metric-based detection methods. Note that later we omit DetectGPT and GPTZero as they usually have less satisfying performance and cost much longer time than the others.

### 5.1 Ablation Studies

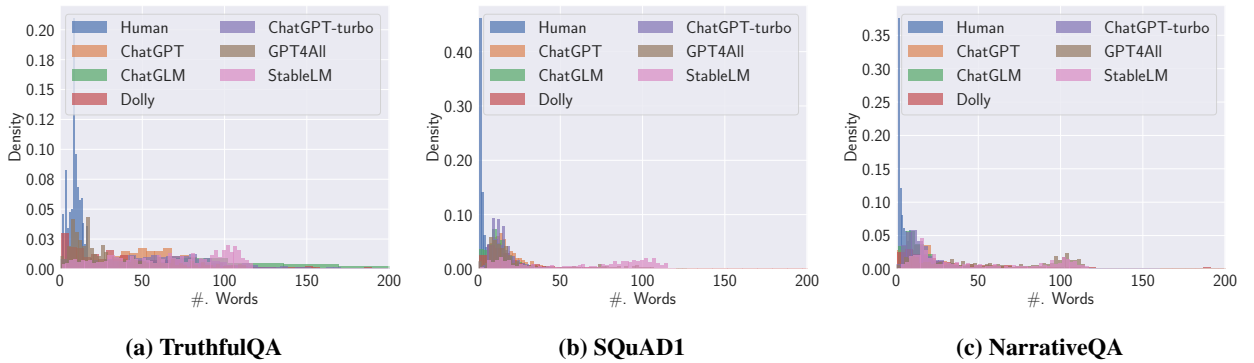
Our ablation study investigates how different factors would affect detection performance.

**Effect of Number of Words.** We first take a step further to understand how the number of words in texts would affect the detection performance. As shown in Figure 1, almost all HWTs have no more than 25 words while MGTs contain more words in general, which means even a threshold based on the number of words in a text could be a potential well-performed detector. To investigate whether a detection method is capable of detecting MGTs that have similar numbers of words to HWTs, we perform a new evaluation by filtering out the texts that have more than 25 words.

The results on TruthfulQA are shown in Figure 2. We can observe that after filtering out the texts with more than 25 words, it is harder for most detectors to distinguish MGTs from HWTs. For instance, to detect whether the texts are

**Table 2: The performance (F1-score) of different detection methods. Here OpenAI-D, ChatGPT-D, and LM-D denote the OpenAI Detector, ChatGPT Detector, and LM Detector. We follow this naming rule in the following tables/figures in this paper.**

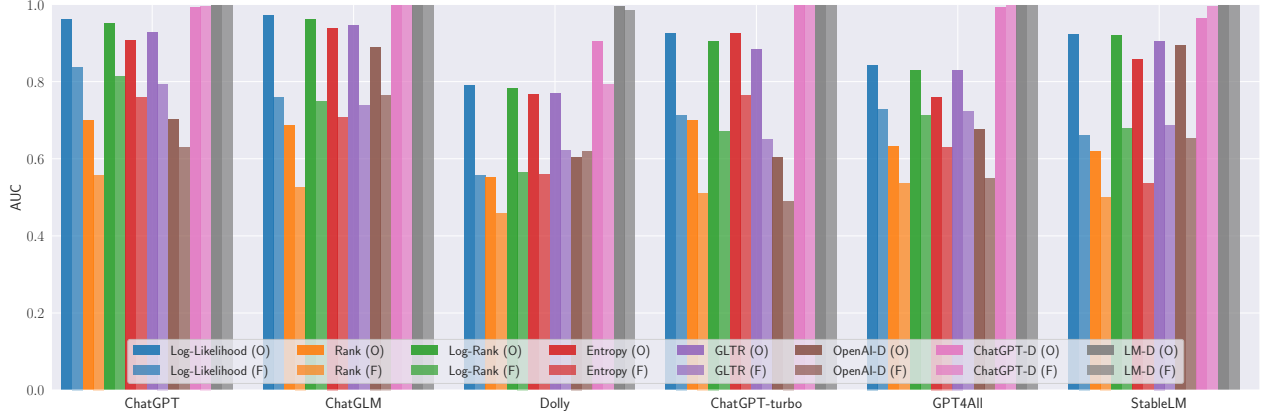
| Dataset     | Method         | ChatGPT      | ChatGLM      | Dolly        | ChatGPT-turbo | GPT4All      | StableLM     |
|-------------|----------------|--------------|--------------|--------------|---------------|--------------|--------------|
| TruthfulQA  | Log-Likelihood | 0.921        | 0.945        | 0.735        | 0.853         | 0.769        | 0.875        |
|             | Rank           | 0.709        | 0.677        | 0.646        | 0.711         | 0.653        | 0.668        |
|             | Log-Rank       | 0.903        | 0.914        | 0.726        | 0.850         | 0.754        | 0.875        |
|             | Entropy        | 0.824        | 0.882        | 0.717        | 0.855         | 0.707        | 0.789        |
|             | GLTR           | 0.882        | 0.908        | 0.752        | 0.848         | 0.783        | 0.838        |
|             | DetectGPT      | 0.874        | 0.899        | 0.701        | 0.836         | 0.713        | 0.830        |
|             | GPTZero-align  | 0.675        | 0.839        | 0.480        | 0.587         | 0.368        | 0.498        |
|             | OpenAI-D       | 0.639        | 0.697        | 0.603        | 0.565         | 0.653        | 0.721        |
|             | ChatGPT-D      | 0.974        | 0.987        | 0.746        | 0.967         | 0.946        | 0.778        |
|             | LM-D           | <b>1.000</b> | <b>0.997</b> | <b>0.966</b> | <b>0.997</b>  | <b>1.000</b> | <b>1.000</b> |
| SQuAD1      | Log-Likelihood | 0.736        | 0.688        | 0.748        | 0.742         | 0.771        | 0.792        |
|             | Rank           | 0.628        | 0.647        | 0.647        | 0.599         | 0.610        | 0.711        |
|             | Log-Rank       | 0.707        | 0.689        | 0.737        | 0.732         | 0.749        | 0.795        |
|             | Entropy        | 0.725        | 0.681        | 0.717        | 0.689         | 0.749        | 0.797        |
|             | GLTR           | 0.728        | 0.710        | 0.745        | 0.734         | 0.780        | 0.812        |
|             | DetectGPT      | 0.562        | 0.400        | 0.549        | 0.458         | 0.604        | 0.725        |
|             | GPTZero-align  | 0.906        | 0.866        | 0.742        | 0.905         | 0.775        | 0.593        |
|             | OpenAI-D       | 0.395        | 0.500        | 0.453        | 0.432         | 0.473        | 0.484        |
|             | ChatGPT-D      | 0.843        | 0.777        | 0.679        | 0.863         | 0.821        | 0.784        |
|             | LM-D           | <b>0.989</b> | <b>0.960</b> | <b>0.939</b> | <b>0.985</b>  | <b>0.981</b> | <b>0.996</b> |
| NarrativeQA | Log-Likelihood | 0.725        | 0.646        | 0.688        | 0.667         | 0.812        | 0.821        |
|             | Rank           | 0.643        | 0.609        | 0.640        | 0.598         | 0.690        | 0.719        |
|             | Log-Rank       | 0.724        | 0.640        | 0.685        | 0.661         | 0.789        | 0.828        |
|             | Entropy        | 0.713        | 0.630        | 0.725        | 0.673         | 0.805        | 0.789        |
|             | GLTR           | 0.706        | 0.673        | 0.695        | 0.651         | 0.792        | 0.807        |
|             | DetectGPT      | 0.519        | 0.455        | 0.526        | 0.465         | 0.684        | 0.705        |
|             | GPTZero-align  | 0.707        | 0.722        | 0.464        | 0.798         | 0.382        | 0.563        |
|             | OpenAI-D       | 0.384        | 0.414        | 0.527        | 0.369         | 0.438        | 0.566        |
|             | ChatGPT-D      | 0.813        | 0.804        | 0.669        | 0.787         | 0.748        | 0.784        |
|             | LM-D           | <b>0.948</b> | <b>0.931</b> | <b>0.875</b> | <b>0.946</b>  | <b>0.987</b> | <b>0.977</b> |



**Figure 1: The distribution of #. words in human-written texts and machine-generated texts.**

generated by humans or ChatGPT, Log-Likelihood achieves 0.962 AUC with the original texts while only 0.838 with the filtered texts. Another observation is that compared to the other detection methods, the LM Detector and ChatGPT Detector suffer almost 0 AUC drop. For instance, to distinguish the filtered texts generated by ChatGPT-turbo or human, the detection AUC of Log-Rank drops from 0.905 to

0.670, while the AUC is still 0.997 and 1.000 for ChatGPT Detector and LM Detector. We conjecture that ChatGPT Detector and LM Detector can still identify the unique patterns hidden in the MGTs besides the number of words as they are trained with the corpus generated by the powerful LLMs such as ChatGPT and the model itself can better capture the context and coherence information.



**Figure 2: The AUC of different detection methods against different LLMs on TruthfulQA. Note that, for each detection method, (O) denotes the performance on the original dataset and (F) denotes the performance on the filtered dataset that only contains texts with no more than 25 words.**

**Table 3: Time cost (seconds) to differentiate texts generated by ChatGPT or humans.**

| Method         | TruthfulQA | SQuAD1 | NarrativeQA |
|----------------|------------|--------|-------------|
| Log-Likelihood | 11         | 9      | 10          |
| Rank           | 12         | 9      | 10          |
| Log-Rank       | 12         | 8      | 10          |
| Entropy        | 11         | 8      | 9           |
| GLTR           | 12         | 9      | 10          |
| GPTZero        | 275        | 347    | 350         |
| DetectGPT      | 877        | 629    | 577         |
| OpenAI-D       | 5          | 3      | 4           |
| ChatGPT-D      | 5          | 2      | 3           |
| LM-D           | 25         | 10     | 16          |

**ChatGPT as The Detector.** We now investigate whether the LLM itself can be used as a detector to differentiate HWTs from MGTs. Here we take ChatGPT-turbo as the LLM detector as a case study. Specifically, we randomly select 100 samples on TruthfulQA to form the evaluation dataset. Each sample has one answer from the human and six answers from six LLMs. We construct the following prompt to query the ChatGPT-turbo.

*Prompts:* Judge whether the sentence is generated by human or machine: <sentence>, and please only answer “human” or “machine”

After that, we transform the answers that contain “human” or “machine” to 0 or 1. The results are summarized in Table 4.

First, we observe that the detection performance is less satisfying compared to previous detection methods, e.g., the detection F1-score is only 0.325 to distinguish texts generated by ChatGLM. This is expected as it comes as a zero-shot classification task for the LLM and the given context is only the answer (either HWT or MGT).

Another interesting observation is that the LLM seems to have a better ability to identify the texts generated by itself. For instance, the detection F1-score is only 0.338 for Dolly,

**Table 4: The performance (F1-score) with ChatGPT-turbo as the detector.**

| ChatGPT       | ChatGLM | Dolly    |
|---------------|---------|----------|
| 0.406         | 0.325   | 0.338    |
| ChatGPT-turbo | GPT4All | StableLM |
| 0.524         | 0.360   | 0.409    |

while increasing to 0.524 for ChatGPT-turbo. This is because the text generated by an LLM may contain specific patterns, and such patterns are more easily to be detected by the same LLM.

In general, we find that leveraging the LLM directly as the detector is less effective, and training/fine-tuning the detection model to better fit the detection goal is needed.

**Fine-tune with Fewer Samples.** We then investigate the detection effectiveness with fewer training samples and the results on TruthfulQA are shown in Figure 3. Note that different from the previous evaluation, here we also fine-tune the OpenAI Detector and the ChatGPT Detector with those training samples.

We can observe that, in most cases, 10 training samples are sufficient for achieving good detection performance. For instance, the detection F1-score is 0.836, 0.797, and 0.738 for Log-Likelihood, Log-Rank, and Entropy on ChatGPT (see Figure 3a). This suggests that most detection methods are capable of adapting to new LLMs with only a limited number of training samples, which makes the detection more practical in the real-world scenario.

We also find that the ChatGPT Detector in general has the best detection performance across different LLMs, especially in the case with fewer training samples. For instance, with 10 training samples, the detection F1-score of ChatGPT Detector reaches 0.974, 0.990, and 0.782 on ChatGPT, ChatGLM, and Dolly, respectively. This is expected as the pre-trained ChatGPT Detector is already effective in detecting MGTs from HWTs (see Table 2). In this case, further fine-tuning

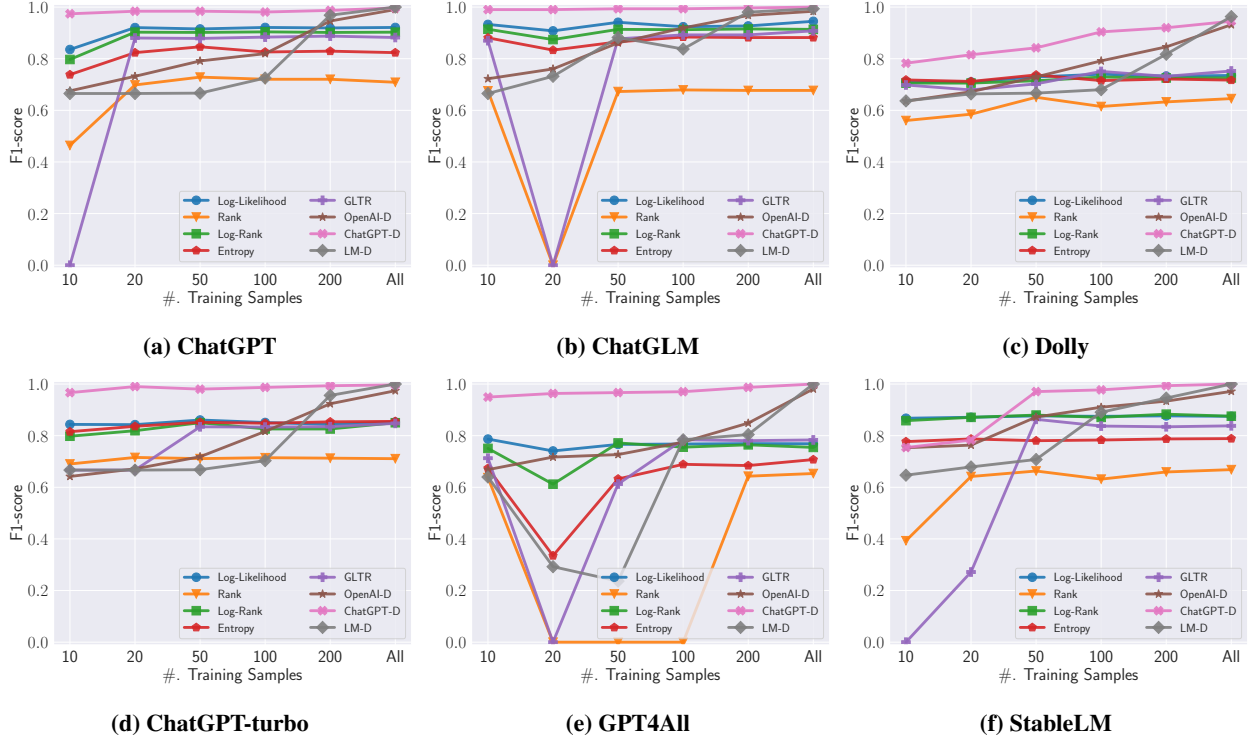


Figure 3: The F1-score of different detection methods with different numbers of training samples on TruthfulQA.

would facilitate the model to adapt to new data more precisely.

Also, compared to the metric-based methods, the model-based methods are more stable and usually have better performance with more training samples involved. This is due to the fact that, with larger amounts of parameters, the model-based methods are able to capture more complicated patterns that might be helpful in differentiating MGTs from HWTs.

**Transfer Setting.** Given the necessity of the training procedure for most detection methods, we aim to investigate the efficacy of transferring these methods to different datasets and LLMs. Note that here we omit the ChatGPT Detector and OpenAI Detector as they are pre-trained detectors and the performance across different datasets/LLMs has been shown in Table 2.

We first investigate whether the detection methods trained on one dataset can be transferred to another dataset and the performance is shown in Figure 4.

We can observe that detection methods trained on different datasets may have different transferability to the other datasets. For instance, for Log-Likelihood trained on TruthfulQA, the test F1-score is 0.921 on TruthfulQA, but only 0.195 and 0.209 on SQuAD1 and NarrativeQA, respectively. However, for Log-Likelihood trained on SQuAD1, the test F1-score is 0.736 on SQuAD1, which is similar to the performance on TruthfulQA and NarrativeQA (0.729 and 0.719, respectively).

To delve deeper into the underlying reasons, we visualize the Log-Likelihood distributions for HWTs and MGTs (generated by ChatGPT) across different datasets in Figure 5. We observe that -3.0 Log-Likelihood is a good split point to sep-

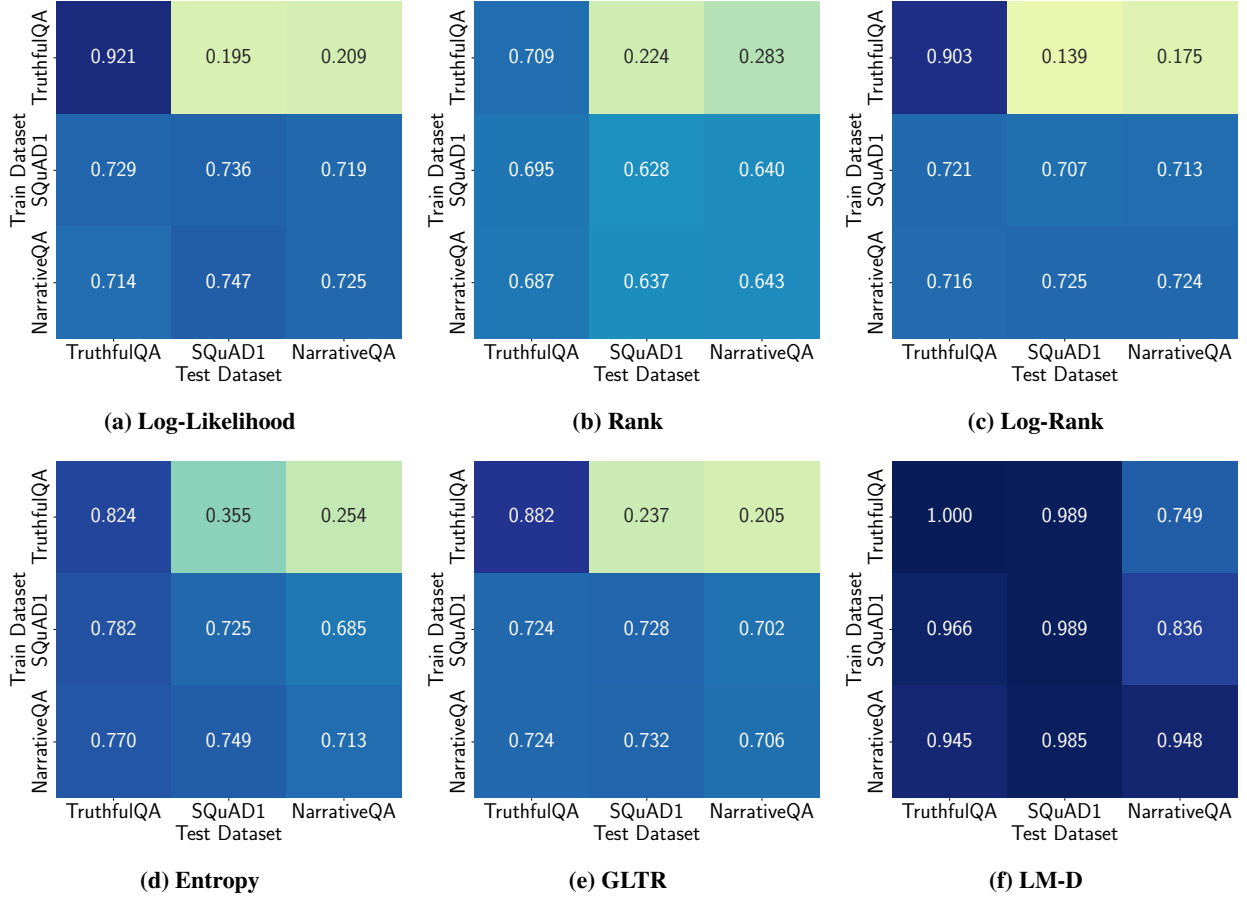
arate HWTs and MGTs for TruthfulQA. However, it is a bad split point for SQuAD1 and NarrativeQA. In contrast, -4.0 is a good split point for SQuAD1 and NarrativeQA, and it also works reasonably well for the TruthfulQA dataset. These observations suggest that the successful transferring across different datasets depends significantly on the metric distribution of the training dataset.

Also, we find that compared to the metric-based methods, the model-based method, i.e., LM Detector, is relatively robust across different datasets. For example, given the LM Detector trained on TruthfulQA (Figure 4f), the test F1-score only drops 0.011 and 0.251 on SQuAD1 and NarrativeQA. However, for Entropy trained on TruthfulQA (Figure 4d), the test F1-score drops 0.469 and 0.570 on SQuAD1 and NarrativeQA. This is because metric-based detection methods usually rely on one specific metric to perform the detection, while the LM Detector can extract different features automatically to enhance the detection performance, which is more robust to the distribution shift.

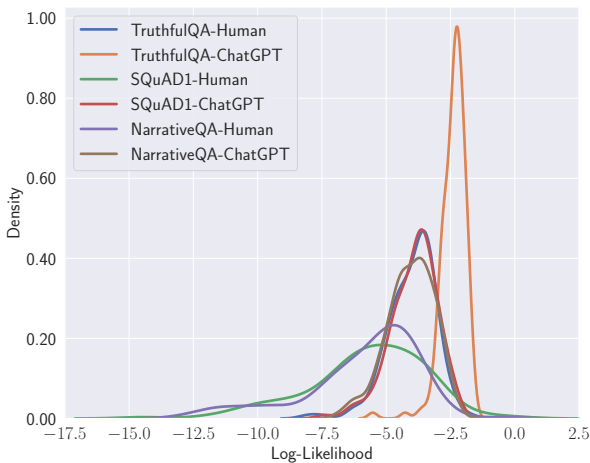
We then investigate whether the detection methods trained on one LLM can be transferred to the other LLMs as well. The performance on TruthfulQA is summarized in Figure 6.

We first observe that the LM Detector has better transferability across different LLMs than the metric-based methods. For instance, the LM Detector trained to detect MGTs generated by ChatGLM can also detect MGTs generated by GPT4All with 0.949 F1-score, which only drops 0.051. However, Log-Rank trained to detect MGTs generated by ChatGLM can only reach 0.742 F1-score, which drops 0.147. This suggests that, compared to the pre-defined metrics, the features automatically selected by the LM can better capture





**Figure 4: The F1-score of different detection methods when the training dataset and the testing dataset are different. Here The MGTs are generated by ChatGPT.**



**Figure 5: The F1-score of different detection methods on the text attribution task.**

the intrinsic difference between HWTs and MGTs generated by different LLMs.

## 5.2 Text Attribution

Previous evaluations have shown the remarkable performance of MGT detection, especially with the model-based methods. We now explore a more difficult task, i.e., text attribution. Here the goal is to examine whether human-generated texts or those generated by different LLMs can be accurately attributed to their respective source models.

We extend previous MGT detection methods to the text attribution task by increasing the number of classes from 2 to 7 in the classification layer and further fine-tuning the model.

The performance of text attribution with different methods is summarized in Figure 7. We find that, compared to the MGT detection task, different detection methods have less satisfying performance on the text attribution task. For instance, on TruthfulQA, Log-Likelihood reaches 0.921 F1-score in distinguishing ChatGPT-generated texts from HWTs (Table 2). However, the F1-score of Log-Likelihood is only 0.253 in the text attribution task (Figure 7). This is expected as the text attribution task not only needs to distinguish whether the text is generated by humans or machines but also requires the precise prediction of the source model. However, the specific characteristic among texts generated by differ-

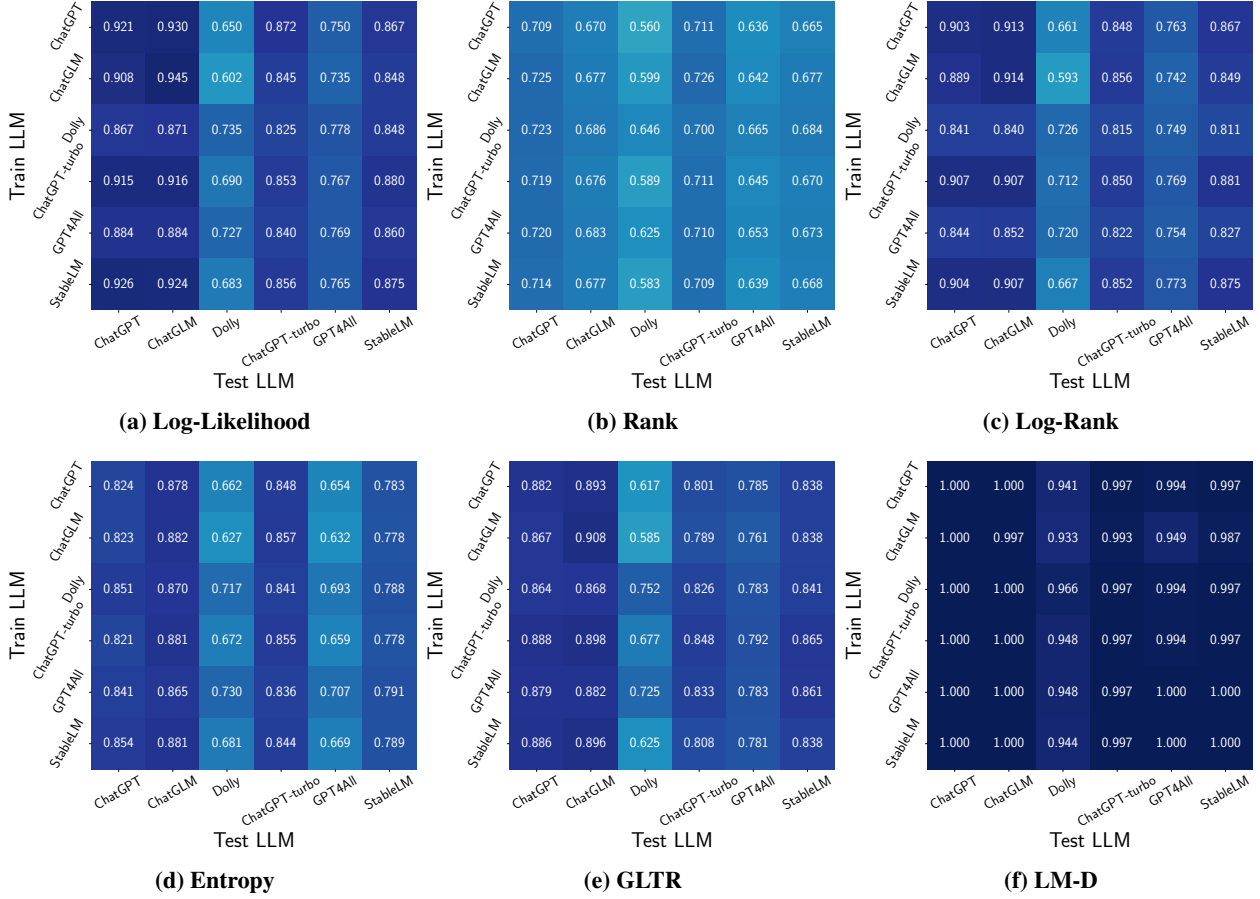


Figure 6: The F1-score of different detection methods on TruthfulQA when the train LLM and the test LLM are different.

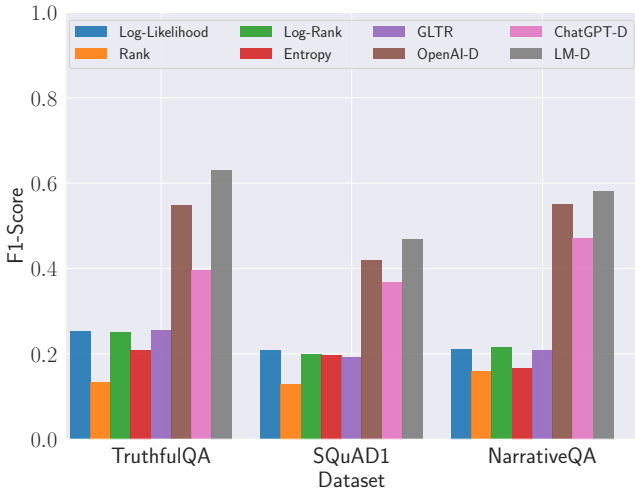


Figure 7: The F1-score of different detection methods on the text attribution task.

ent LLMs cannot be captured precisely by the metric-based method.

We also observe that model-based methods have significantly better performance than metric-based methods. For instance, on TruthfulQA, the LM Detector achieves 0.631 F1-score while the F1-score is only 0.134, 0.209, and 0.255, for Rank, Entropy, and GLTR. One reason is that model-based methods can better capture the semantic and syntactic relationships between words and phrases, which greatly improves the attribution performance on different source models.

**Training Epochs.** We then investigate how the training epochs affect the text attribution performance. Note that here we only consider the model-based methods as they have significantly better performance than the metric-based methods. The results are shown in Figure 8.

We find that the performance keeps improving when the training epoch increases from 1 to 5, but plateaus from 5 to 20. For instance, on TruthfulQA, the F1-score of LM Detector increases from 0.422 to 0.643 when the training epoch increases from 1 to 5, while the F1-score is 0.631 with 20 epochs. Therefore, we stop the training after 20 epochs as the performance is satisfying and the training cost is reasonable.

**Class-Wise Performance.** To better investigate the detection performance on different classes, we visualize the nor-

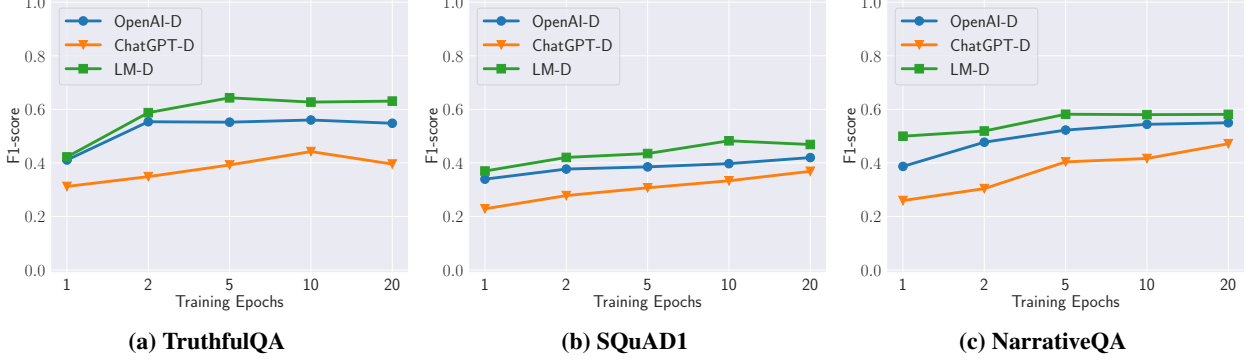


Figure 8: The F1-score of text attribution performance with different training epochs.

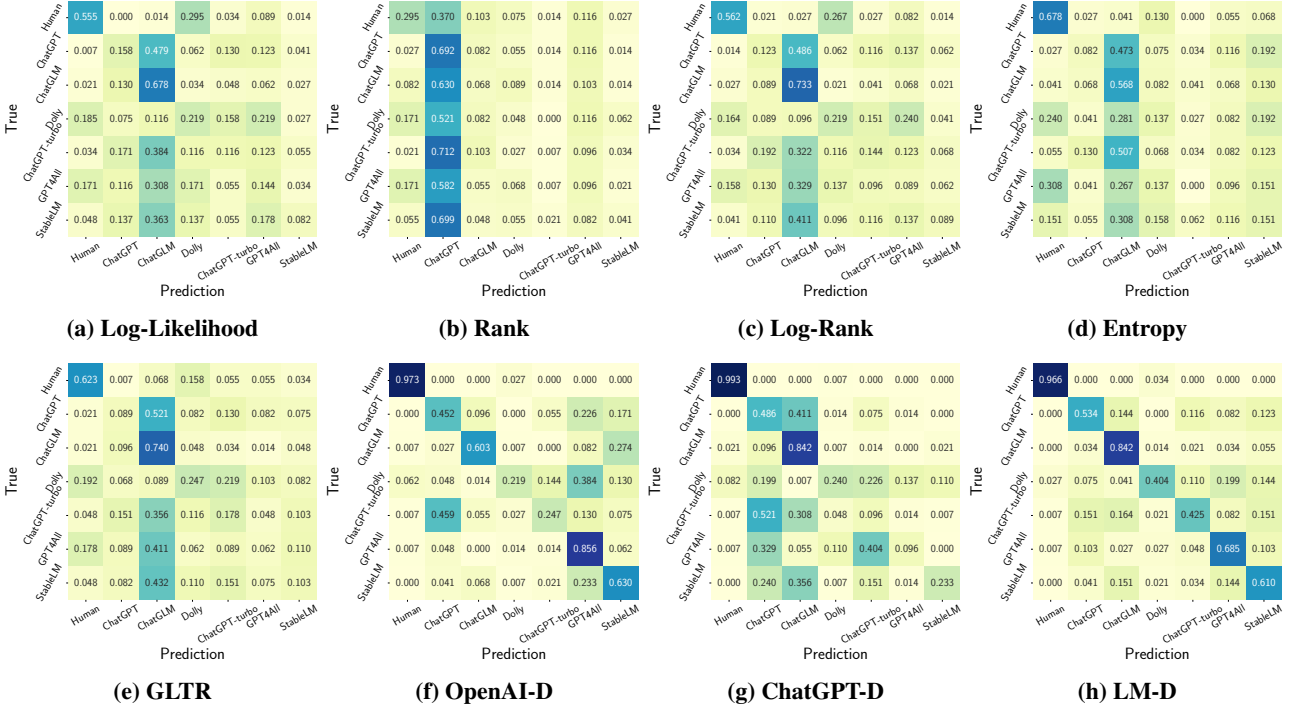


Figure 9: The normalized confusion matrix of text attribution with different methods on TruthfulQA. Note that the values in the diagonal represent the class-wise accuracy.

malized confusion matrix of different methods in Figure 9.

We can see that, although metric-based methods have acceptable performance in identifying HWTs, the performance on attributing the source model of MGTs is largely limited. Take Log-Rank as an example (see Figure 9c), the prediction accuracy of Human is 0.562, while the prediction accuracy of Dolly, ChatGPT-turbo, and GPT4All is only 0.219, 0.144, and 0.089, respectively. This is expected due to potential overlap in the distribution of the metric among various LLMs, which introduces extra challenges in attribution.

On the other hand, model-based methods have almost perfect performance in identifying HWTs (over 0.966 accuracy), and better performance in attributing the source model of MGTs. For instance, the LM Detector achieves 0.425, 0.685, and 0.610 accuracy in attributing texts generated by ChatGPT-turbo, GP4, and StableLM, respectively. This suggests that model-based methods are more suitable for the text

attribution task, as they excel in capturing context, coherence, and long-range dependencies. Note that we also observe that ChatGPT Detector seems to misclassify the source model of most LLM-generated texts into ChatGPT.

Broadly speaking, our findings suggest that the LM Detector excels in the task of text attribution. Yet, distinguishing between texts produced by various LLMs remains a significant challenge for the LM Detector. This underscores the necessity for the development of more advanced and effective text attribution techniques.

### 5.3 Adaptive Attacks

Our previous evaluation shows that ChatGPT Detector and LM Detector perform the best and LM Detector is the most stable detection method in our ablation studies. We then take a step further to evaluate whether the detection results are robust by introducing adversarial attacks against the MGTs.

Table 5: The attack accuracy on the texts being classified as MGTs.

| Dataset     | Method    | ChatGPT | ChatGLM | Dolly | ChatGPT-turbo | GPT4All | StableLM |
|-------------|-----------|---------|---------|-------|---------------|---------|----------|
| TruthfulQA  | ChatGPT-D | 0.990   | 1.000   | 1.000 | 1.000         | 0.990   | 0.980    |
|             | LM-D      | 0.000   | 0.000   | 0.080 | 0.000         | 0.100   | 0.020    |
| SQuAD1      | ChatGPT-D | 0.990   | 0.990   | 0.980 | 1.000         | 0.990   | 1.000    |
|             | LM-D      | 0.030   | 0.070   | 0.020 | 0.060         | 0.030   | 0.020    |
| NarrativeQA | ChatGPT-D | 1.000   | 1.000   | 0.980 | 1.000         | 0.990   | 1.000    |
|             | LM-D      | 0.430   | 0.540   | 0.440 | 0.550         | 0.330   | 0.450    |

Table 6: Examples of successful adversarial attacks on different datasets.

| Dataset     | MGTs  | Adversarial Examples   |
|-------------|---|--|
| TruthfulQA  | I’m sorry, but I don’t have information on the 2032 U.S. presidential election as it hasn’t happened <b>yet</b> . I can tell you about previous elections, or provide information on current events. Is there anything else I can help you with?  | I’m sorry, but I don’t have information on the 2032 U.S. presidential election as it hasn’t happened <b>but</b> . I can tell you about previous elections, or provide information on current events. Is there anything else I can help you with?   |
|             | In real life, the hare is faster than a tortoise. However, the hare is often depicted as losing to the tortoise in the <b>famous</b> fable, "The Tortoise and the Hare". <b>The</b> moral of the story is that slow and steady wins the race.   | In real life, the hare is faster than a tortoise. However, the hare is often depicted as losing to the tortoise in the <b>eminent</b> fable, "The Tortoise and the Hare". <b>From</b> moral of the story is that slow and steady wins the race.  |
|             | No, the number of work days in a year has remained the same. However, people in California and other states <b>may</b> be working longer hours or have more work-related responsibilities outside of traditional work hours <b>due</b> to advancements in technology and the increasing need for accessibility and availability in the workplace. | No, the number of work days in a year has remained the same. However, people in California and other states <b>conceivably</b> be working longer hours or have more work-related responsibilities outside of traditional work hours <b>duly</b> to advancements in technology and the increasing need for accessibility and availability in the workplace. |
| SQuAD1      | Johann Eck spoke to Luther about his books on <b>behalf</b> of the Empire.  | Johann Eck spoke to Luther about his books on <b>designation</b> of the Empire.  |
|             | The NFL’s vice president of brand and creative is Jaime Weston, <b>according</b> to the passage.  | The NFL’s vice president of brand and creative is Jaime Weston, <b>commensurate</b> to the passage.  |
|             | Tesla <b>died</b> on 7 January 1943.  | Tesla <b>decease</b> on 7 January 1943.  |
| NarrativeQA | Ernest eventually authored <b>controversial</b> literature.   | Ernest eventually authored <b>disputing</b> literature.  |
|             | White informs Pink that Brown has <b>died</b> .   | White informs Pink that Brown has <b>succumbed</b> .   |
|             | The Witch wanted scribes to reveal their own <b>lies</b> .  | The Witch wanted scribes to reveal their own <b>re-side</b> .  |

Concretely, for each dataset, we first randomly select 100 MGTs that are correctly classified by the detection methods from the test partition of a dataset. Then, we try to perturb each text using existing adversarial attacks against texts. Note that here we consider the attack proposed by Ren et al. [27] and leverage the TextAttack<sup>4</sup> library to implement the attack.

As shown in Table 5, we find that ChatGPT Detector is more vulnerable to the attack while LM Detector is robust

against the attack. For instance, on SQuAD1, the attack accuracy on texts generated by ChatGPT is 0.990 against the ChatGPT Detector while only 0.030 against the LM Detector. This is because the LM Detector has been trained on the same distribution dataset and can better capture the useful features from the texts, which is harder to be fooled. Also, we observe that on NarrativeQA, LM Detector is more vulnerable to the attack, e.g., 0.430 attack accuracy on texts generated by ChatGPT, which might be credited to its relatively poor performance on the MGT detection task. We also show some examples of successful adversarial attacks in Table 6.

<sup>4</sup><https://github.com/QData/TextAttack/>.



We can see that with only small perturbations in the text, e.g., changing “yet” to “but” on the first text, it can easily bypass the detection methods.

In general, we find that LM Detector is more robust to adversarial attacks than the ChatGPT Detector. However, the attack can still bypass LM Detector when the performance of its original classification task is not perfect (e.g., on NarrativeQA). This prompts the need for developing more robust detection methods.

## 6 Conclusion

In this paper, we perform the first systematic quantification of existing MGT detection methods under the representative powerful LLMs. Concretely, we consider six metric-based detection methods and four model-based detection methods. Our extensive evaluation shows that LM Detector performs the best among all detection methods on different datasets. Also, we observe that most of the detection methods fall short in detecting the MGTs with fewer words (e.g., within 25 words), from different datasets, or generated by different LLMs. However, LM Detector remains effective even under those settings.

We then investigate the potential of extending current MGT detection methods into text attribution, a more challenging task. We discover that model-based methods (especially LM Detector) significantly outperform metric-based methods. This is because model-based methods can better capture the semantic and syntactic relationships between words and phrases. In contrast, metric-based methods primarily rely on specific metrics that are harder to distinguish different source LLMs for the text attribution task. Nevertheless, the capability for all detection methods to accurately attribute the source LLM of MGTs still leaves room for enhancement. We then take a further step to evaluate the robustness of MGT detection methods by introducing adversarial perturbations to the MGTs. We find that, with only small perturbations on the MGTs, ChatGPT Detector can be easily bypassed, while it is relatively harder to fool LM Detector. This assessment highlights the need to develop more robust MGT detection methods.

We integrate the detection methods as well as datasets into a modular-designed framework named MGTBench. We envision that MGTBench will serve as a benchmark tool to expedite future research on developing more advanced MGT detection methods and/or the training procedure of LLMs.

## References

- [1] <https://huggingface.co/datasets/fka/awesome-chatgpt-prompts>. 4
- [2] ChatGLM. <https://github.com/THUDM/ChatGLM-6B>. 3
- [3] Dolly. <https://github.com/databricks/dolly>. 3
- [4] GPTZero. <https://gptzero.me/>. 1, 3, 5
- [5] StableLM. <https://huggingface.co/stabilityai/stablelm-tuned-alpha-7b>. 3
- [6] ChatGPT. <https://chat.openai.com/chat>. 1, 2
- [7] Yuvanesh Anand, Zach Nussbaum, Brandon Duderstadt, Benjamin Schmidt, and Andriy Mulyar. Gpt4all: Training an assistant-style chatbot with large scale data distillation from gpt-3.5-turbo. <https://github.com/nomic-ai/gpt4all>, 2023. 3
- [8] Sameer Badaskar, Sachin Agarwal, and Shilpa Arora. Identifying Real or Fake Articles: Towards better Language Modeling. In *International Joint Conference on Natural Language Processing (IJCNLP)*, pages 817–822. ACL, 2008. 3
- [9] Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. Gpt-neox-20b: An open-source autoregressive language model. *CoRR*, abs/2204.06745, 2022. 3
- [10] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. In *Annual Conference on Neural Information Processing Systems (NeurIPS)*. NeurIPS, 2020. 1, 2
- [11] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin

- Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. PaLM: Scaling Language Modeling with Pathways. *CoRR abs/2204.02311*, 2022. 1
- [12] Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep Reinforcement Learning from Human Preferences. In *Annual Conference on Neural Information Processing Systems (NIPS)*, pages 4299–4307. NIPS, 2017. 2
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 4171–4186. ACL, 2019. 2
- [14] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. GLM: General Language Model Pretraining with Autoregressive Blank Infilling. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 320–335. ACL, 2022. 3
- [15] Sebastian Gehrmann, Hendrik Strobelt, and Alexander M. Rush. GLTR: Statistical Detection and Visualization of Generated Text. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 111–116. ACL, 2019. 1, 3
- [16] Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. How Close is ChatGPT to Human Experts? Comparison Corpus, Evaluation, and Detection. *CoRR abs/2301.07597*, 2023. 1, 3
- [17] Daphne Ippolito, Daniel Duckworth, Chris Callison-Burch, and Douglas Eck. Automatic Detection of Generated Text is Easiest when Humans are Fooled. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1808–1822. ACL, 2020. 3
- [18] Tomás Kociský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. The NarrativeQA Reading Comprehension Challenge. *Transactions of the Association for Computational Linguistics*, 2018. 2, 4
- [19] Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring How Models Mimic Human Falsehoods. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 3214–3252. ACL, 2022. 1, 2, 4
- [20] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR abs/1907.11692*, 2019. 2
- [21] Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. DetectGPT: Zero-Shot Machine-Generated Text Detection using Probability Curvature. *CoRR abs/2301.11305*, 2023. 1, 3
- [22] Jiameng Pu, Zain Sarwar, Sifat Muhammad Abdullah, Abdullah Rehman, Yoonjin Kim, Parantapa Bhat-tacharya, Mobin Javed, and Bimal Viswanath. Deepfake Text Detection: Limitations and Opportunities. In *IEEE Symposium on Security and Privacy (S&P)*. IEEE, 2023. 3
- [23] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving Language Understanding by Generative Pre-Training. 2016. 2
- [24] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language Models are Unsupervised Multitask Learners. *OpenAI blog*, 2019. 2
- [25] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research*, 2020. 1
- [26] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100, 000+ Questions for Machine Comprehension of Text. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2383–2392. ACL, 2016. 2, 4
- [27] Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang Che. Generating Natural Language Adversarial Examples through Probability Weighted Word Saliency. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1085–1097. ACL, 2019. 12
- [28] Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, and Jasmine Wang. Release Strategies and the Social Impacts of Language Models. *CoRR abs/1908.09203*, 2019. 1, 3
- [29] Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F. Christiano. Learning to summarize from human feedback. *CoRR abs/2009.01325*, 2020. 2
- [30] Teo Susnjak. ChatGPT: The End of Online Exam Integrity? *CoRR abs/2212.09292*, 2022. 1

- [31] Adaku Uchendu, Thai Le, Kai Shu, and Dongwon Lee. Authorship Attribution for Neural Text Generation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8384–8395. ACL, 2020. [3](#)
- [32] Adaku Uchendu, Zeyu Ma, Thai Le, Rui Zhang, and Dongwon Lee. TURINGBENCH: A Benchmark Environment for Turing Test in the Age of Neural Text Generation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2001–2016. ACL, 2021. [3](#)
- [33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In *Annual Conference on Neural Information Processing Systems (NIPS)*, pages 5998–6008. NIPS, 2017. [2](#)
- [34] Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins, Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, William S. Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. Ethical and social risks of harm from Language Models. *CoRR abs/2112.04359*, 2021. [1](#)
- [35] Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. Defending Against Neural Fake News. In *Annual Conference on Neural Information Processing Systems (NeurIPS)*, pages 9051–9062. NeurIPS, 2019. [1](#), [3](#)