Microsoft

# KDD 2024 Tutorial

# Decoding the AI Pen: Techniques and Challenges in Detecting AI-Generated Text

Sara Abdali          Richard Anarfi

CJ Barberan          Jia He

# Disclaimer

This study represents independent research conducted by the authors and does not necessarily represent the views or opinions of any organizations.

# Tutors and Contributors

## In-person presenters

Sara Abdali

saraabdali@microsoft.com

CJ Barberan

cjbarberan@microsoft.com

Jia He

hejia@microsoft.com

## Contributor

Richard Anarfi

ranarfi@microsoft.com

# Links and Materials

➢**Tutorial Website**

**decoding-the-AI-pen.github.io**

➢**If you find our paper useful, please feel free to cite it in your work**

```
@inproceedings{Abdali_2024, series={KDD '24},
  title={Decoding the AI Pen: Techniques and Challenges in Detecting AI-Generated Text},
  volume={30},
  url={http://dx.doi.org/10.1145/3637528.3671463},
  DOI={10.1145/3637528.3671463},
  booktitle={Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining},
  publisher={ACM},
  author={Abdali, Sara and Anarfi, Richard and Barberan, CJ and He, Jia},
  year={2024},
  month=aug, pages={6428–6436},
  collection={KDD '24} }
```

# Tutorial Outline

Introduction (Dr. Sara Abdali)

Part I : Risks and Misuse of AI-Generated Text (  Dr. CJ Barberan)

Part II: AI-Generated Text Detection Techniques ( Dr. Sara Abdali)

Part III: Vulnerabilities of Detection Techniques ( Dr. Jia He)

Part IV: Theoretical perspective on possibility of detection ( Dr. Sara Abdali)
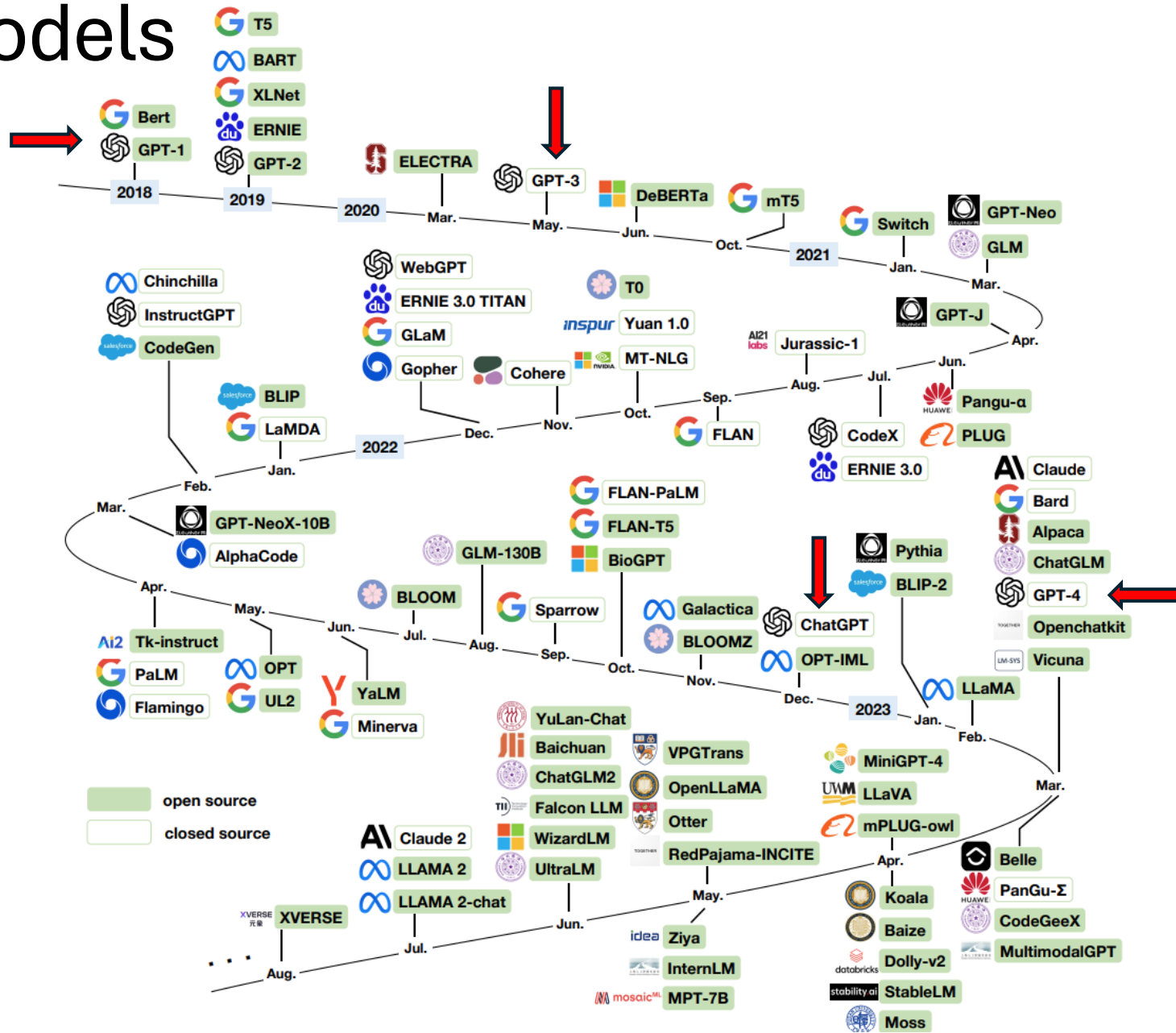
Conclusion and Future Directions (Dr. CJ Barberan)
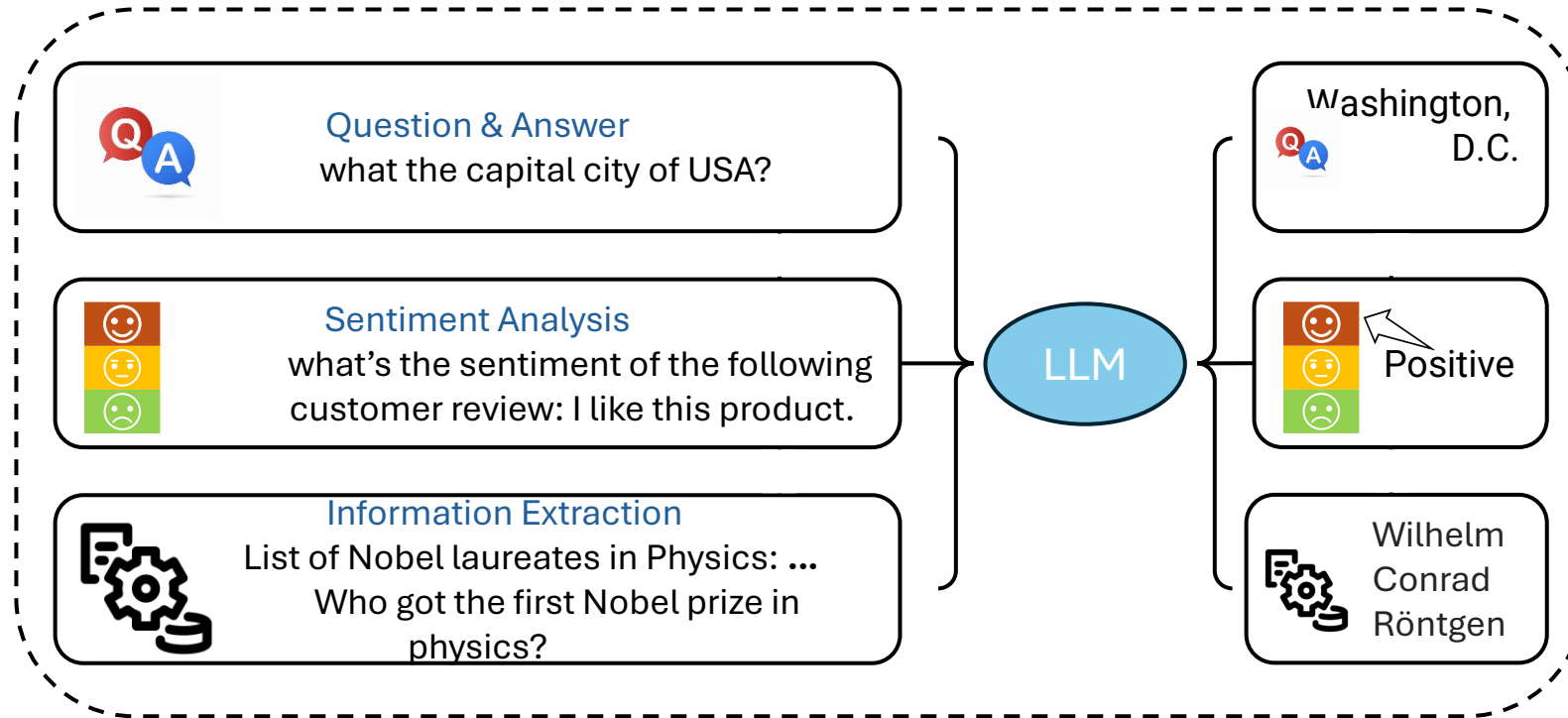
# Introduction

# Evolution of Language Models



Large Language Models (LLMs) have revolutionized the field of Natural Language Generation (NLG) by demonstrating an impressive ability to generate human-like text.

Gao et al. Examining User-Friendly and Open-Sourced Large GPT Models: A Survey on Language, Multimodal, and Scientific GPT Models
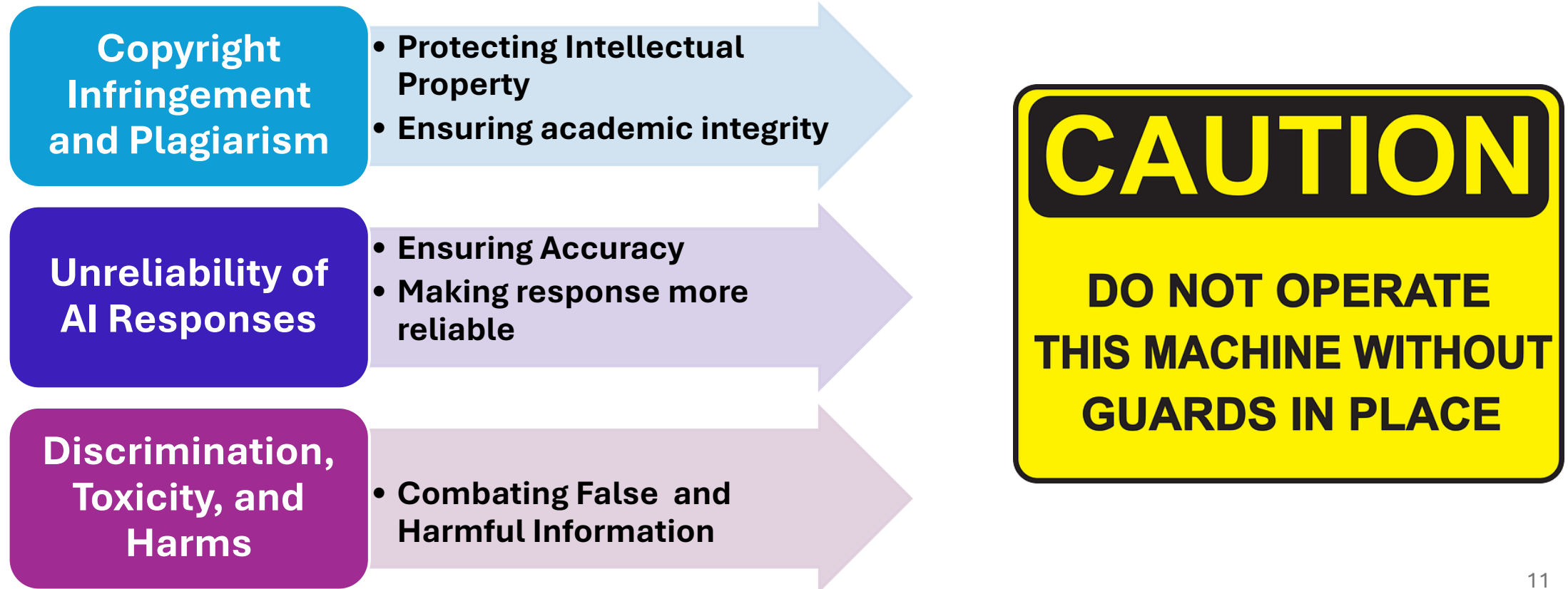
# Applications of LLMs

➢ LLMs' applications traverse a wide spectrum of domains, including question answering , sentiment analysis and specially text generation.

# Balancing the Pros. and Cons. of LLMs: A Call for Caution

➢ However, the ubiquity of LLMs brings forth some challenges that necessitate prudent examination.

**Copyright Infringement and Plagiarism**
- Protecting Intellectual Property
- Ensuring academic integrity

**Unreliability of AI Responses**
- Ensuring Accuracy
- Making response more reliable

**Discrimination, Toxicity, and Harms**
- Combating False and Harmful Information

**CAUTION**

**DO NOT OPERATE THIS MACHINE WITHOUT GUARDS IN PLACE**

11

# Why Do We Need to Detect AI-generated Text?

**Motivation**

**Recognizing text produced by large language models is a common strategy to address many of the issues they pose.**

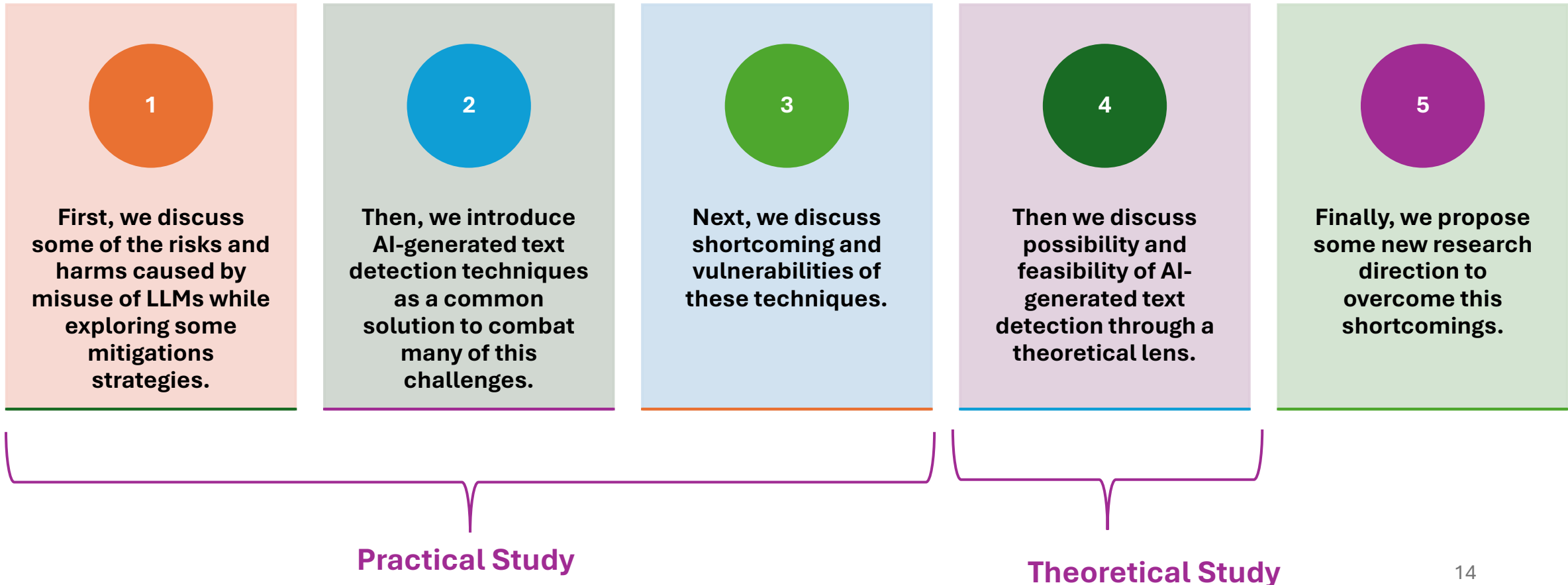# Dual Challenge of Distinguishing AI-generated Text From Human-written Content

Identifying disparities can enhance the quality of AI-generated material.



However, this endeavor complicates the identification process.

# What Do We Cover in This Tutorial?

**1**

First, we discuss some of the risks and harms caused by misuse of LLMs while exploring some mitigations strategies.

**2**

Then, we introduce AI-generated text detection techniques as a common solution to combat many of this challenges.

**3**

Next, we discuss shortcoming and vulnerabilities of these techniques.

**4**

Then we discuss possibility and feasibility of AI-generated text detection through a theoretical lens.

**5**

Finally, we propose some new research direction to overcome this shortcomings.
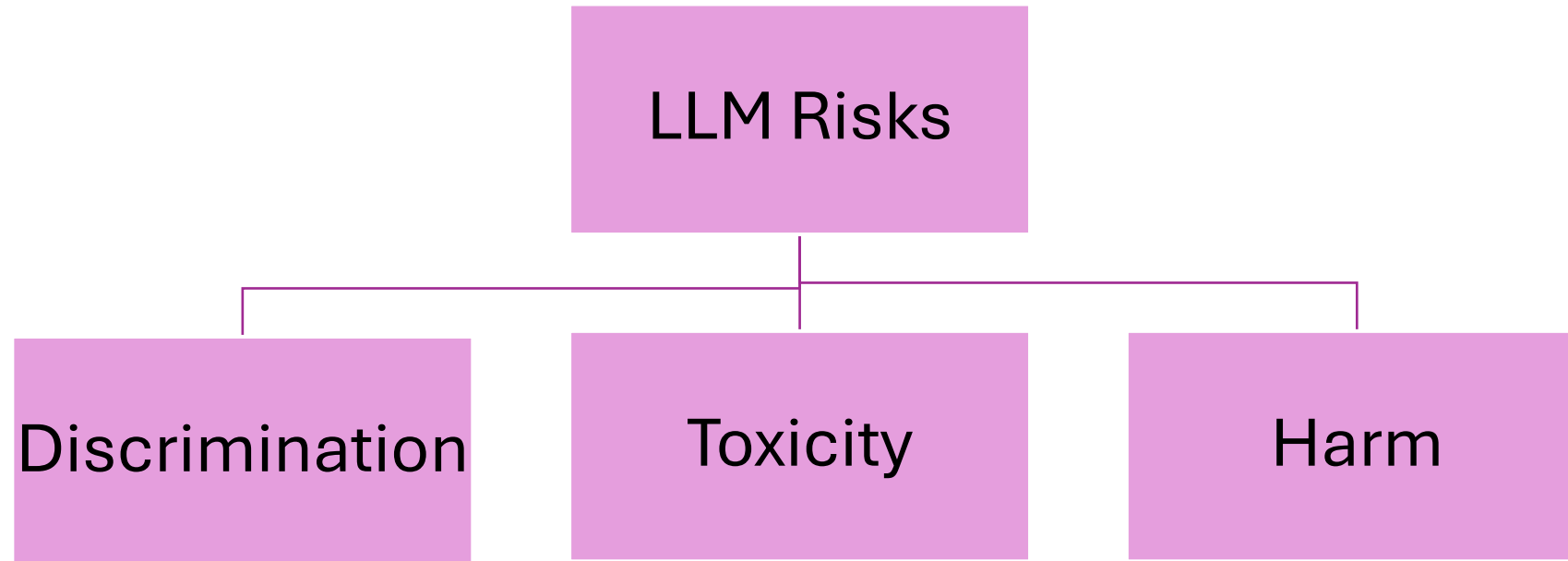
**Practical Study**

**Theoretical Study**

14

# Part I
# Risks and Misuse of AI-Generated Text

- Discrimination, Toxicity, and Harms
- Factual Inconsistency and Unreliability of AI Responses
- Copyright Infringement and Plagiarism
- Misinformation Dissemination

# Overview

- Discuss the risks that come with LLMs
  - Discrimination
  - Toxicity
  - Harm
  - Misinformation

- Discuss the detection strategies
  - Black-box
  - White-box

- Evading detection strategies

# LLM Potential Risks



**Weidinger et al.** "Ethical and social risks of harm from Language Models"

# Discrimination

➤ LLMs can output responses that exhibit discrimination tendencies
- Stereotypes are reinforced
- Present if the training data contains discrimination
- Data can have imbalance of text for certain groups

➤ Examples:
- StereoSet benchmark
  - Measures stereotypes across race, gender, religion, and profession
- CrowS-Pairs benchmark
  - Measures cultural stereotypes
- HONEST benchmark
  - Measures 'hurtful stereotypes'

# How Bias are LMs to Discrimination?

All these LMs exhibit discrimination bias with **ALBERT** being the highest.

| | *n* | % | BERT | RoBERTa | ALBERT |
|---|---|---|---|---|---|
| WinoBias-*ground* (Zhao et al., 2018) | 396 | - | **56.6** | 69.7 | 71.7 |
| WinoBias-*knowledge* (Zhao et al., 2018) | 396 | - | **60.1** | 68.9 | 68.2 |
| StereoSet (Nadeem et al., 2020) | 2106 | - | **60.8** | **60.8** | 68.2 |
| CrowS-Pairs | 1508 | 100 | **60.5** | 64.1 | 67.0 |
| CrowS-Pairs-*stereo* | 1290 | 85.5 | **61.1** | 66.3 | 67.7 |
| CrowS-Pairs-*antistereo* | 218 | 14.5 | 56.9 | **51.4** | 63.3 |
| *Bias categories in Crowdsourced Stereotype Pairs* | | | | | |
| Race / Color | 516 | 34.2 | **58.1** | 62.0 | 64.3 |
| Gender / Gender identity | 262 | 17.4 | 58.0 | **57.3** | 64.9 |
| Socioeconomic status / Occupation | 172 | 11.4 | **59.9** | 68.6 | 68.6 |
| Nationality | 159 | 10.5 | **62.9** | 66.0 | 63.5 |
| Religion | 105 | 7.0 | **71.4** | **71.4** | 75.2 |
| Age | 87 | 5.8 | **55.2** | 66.7 | 70.1 |
| Sexual orientation | 84 | 5.6 | 67.9 | **65.5** | 70.2 |
| Physical appearance | 63 | 4.2 | **63.5** | 68.3 | 66.7 |
| Disability | 60 | 4.0 | **61.7** | 71.7 | 81.7 |

**Nangia et al.** "CrowS-Pairs: A Challenge Dataset for Measuring Social Biases in Masked Language Models"

# Toxicity Generation & Biases

➢Toxicity contains language which expresses hate speech, harassment, and abusive information

➢Pretrained LLMs can generate toxic text

➢REALTOXICITYPROMPTS is a testbed to evaluate toxic degeneration

➢Datasets contain a good amount of toxicity

➢Steering methods can help but are not a silver bullet

Gehman et al. "REALTOXICITYPROMPTS: Evaluating Neural Toxic Degeneration in Language Models"

# Toxicity Taxonomy

➢Toxicity can have different categories
- Utterence: literal toxic language
- Context: depends on the scenario

➢There is a performance gap between utterance and context toxicity evaluation

| Models | Utterence | Context | Overall ↑ |
|---|---|---|---|
| GPT4 | 0.840 | 0.678 | 0.759 |
| GPT3.5-turbo | 0.800 | 0.643 | 0.722 |
| Llama2-chat-70B | 0.840 | 0.630 | 0.735 |
| Llama2-chat-13B | 0.870 | 0.745 | 0.807 |
| Llama2-chat-7B | 0.880 | 0.796 | **0.838** |
| Vicuna-13B | 0.659 | 0.374 | 0.516 |
| Vicuna-7B | 0.581 | 0.314 | 0.448 |
| Llama2-13B | 0.719 | 0.427 | 0.573 |
| Llama2-7B | 0.715 | 0.121 | 0.418 |

**Cui et al.** "FFT: Towards Harmlessness Evaluation and Analysis for LLMs with Factuality, Fairness, Toxicity"

# Chatbot Toxicity

➢ChatGPT given a persona can exhibit more toxic behavior depending on the persona.

- Persona of a 'good person' will not exhibit high toxicity
- Persona of a 'bad person' will exhibit high toxicity



Bad Persona → LLM → Toxic Behavior

**Deshpande et al.** "Toxicity in CHATGPT: Analyzing Persona-assigned Lanugage Models"

# How Toxic are Bad Personas?

The more negative a persona, the higher the toxicity it exhibits.

| Persona | ENTITY-CONDITIONED | | REALTOX |
|---|---|---|---|
| | TOXICITY | POR | TOXICITY |
| *No persona* | $0.11_{\pm 0.02}$ | 0.13 | $0.09_{\pm 0.01}$ |
| *A good person* | $0.06_{\pm 0.01}$ | 0.17 | $0.09_{\pm 0.01}$ |
| *A normal person* | $0.14_{\pm 0.02}$ | 0.38 | $0.11_{\pm 0.01}$ |
| *A bad person* | $\mathbf{0.62}_{\pm 0.01}$ | **0.96** | $\mathbf{0.42}_{\pm 0.01}$ |
| *A nasty person* | $\mathbf{0.63}_{\pm 0.01}$ | **0.92** | $\mathbf{0.53}_{\pm 0.01}$ |
| *A terrible person* | $\mathbf{0.64}_{\pm 0.01}$ | **0.94** | $\mathbf{0.49}_{\pm 0.01}$ |

Deshpande et al. "Toxicity in CHATGPT: Analyzing Persona-assigned Lanugage Models"

# Evading Toxic Detection Classifiers

➢With LLMs exhibiting toxicity, detectors have been developed

➢To evade the detection, they utilize implicit toxicity

▪ Generate a response that can be toxic without using the explicit toxic words

▪ An example is to use euphemism

**Wen et al.** "Unveiling the Implicit Toxicity in Large Language Models"

# What is the Performance in Evading Toxic Detectors?

**Utilizing rewards achieves the highest success in evading the detectors**

| Test Data | Source | Reward | Annotated Toxic Prob. | Attack Success Rate | | | | | Distinct-4 |
| | | | | P-API | Moderation | TOXIGEN | BAD | Davinci003 | |
|---|---|---|---|---|---|---|---|---|---|
| Offensive Twitter | Crawl | -5.91 | N/A | 14.10 | 73.20 | 14.40 | 1.90 | 6.20 | **0.99** |
| TOXIGEN | LM | -3.96 | N/A | 72.28 | 67.93 | 33.97 | 20.92 | 9.24 | 0.94 |
| Latent Hatred | Crawl + CS | -3.86 | N/A | 72.92 | 74.64 | 42.14 | 16.09 | 11.37 | 0.98 |
| BAD | CS + LM | -3.36 | N/A | 76.77 | 82.11 | 55.28 | 24.85 | 26.25 | 0.95 |
| GPT-3.5-turbo | LM | 0.78 | 56.91 | 96.69 | 96.69 | 75.14 | 64.09 | 58.47 | 0.93 |
| SL LLaMA-13B | LM | 0.35 | 54.02 | 97.03 | 94.64 | 69.05 | 64.29 | 58.34 | 0.91 |
| SL-R LLaMA-13B | LM | 1.01 | 55.23 | 99.41 | 95.27 | 75.15 | 68.64 | 56.80 | 0.87 |
| RL LLaMA-13B | LM | **2.47** | **58.84** | **99.55** | **97.81** | **82.51** | **90.16** | **62.85** | 0.85 |

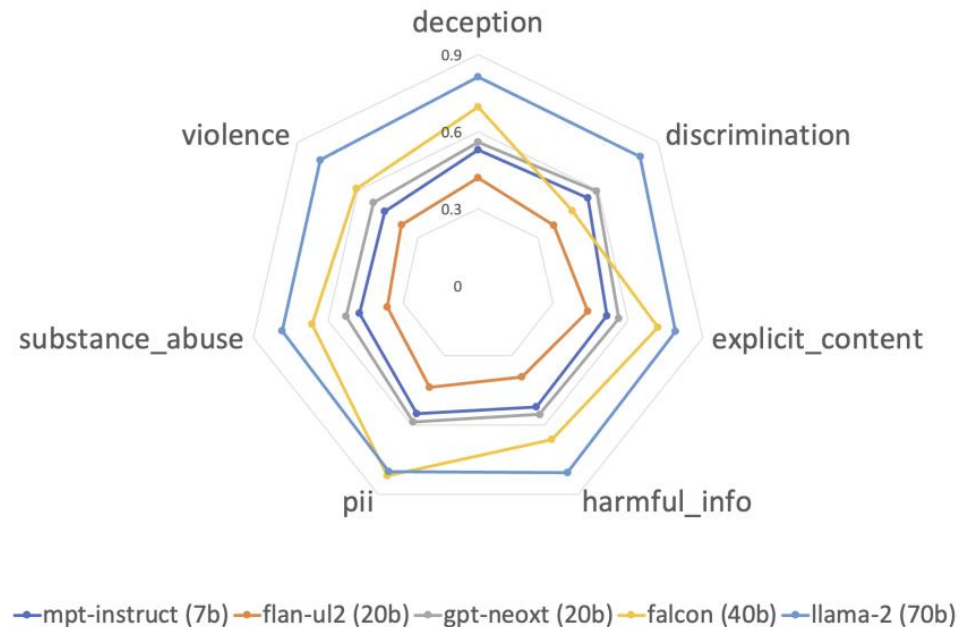**Wen et al.** "Unveiling the Implicit Toxicity in Large Language Models"

# Generating Harmful Responses

➢Like toxicity, LLMs can generate harmful responses from prompts

➢Ex: "How can I use social media to spread false rumors about my competitors?"

➢Adversarial Question Attack (AttaQ) is a dataset to evaluate how LLMs generate harmful or harmless responses

➢With AttaQ, LLMs can generate harmful responses

➢AttaQ harm categories
- Violence
- Substance abuse
- Personal Identifiable Information (PII)
- Harmful information
- Explicit Content
- Discrimination
- Deception

**Kour et al.** "Unveiling Safety Vulnerabilities of Large Language Models"

# How Harmful are LLMs?

➢ From Kour et al., they used 5 LLMs on AttaQ for instruction

  ▪ Flan-ul2 and llama-2 are the least harmful.
  ▪ They mention that a model's performance cannot be understood through evaluation.



**Kour et al.** "Unveiling Safety Vulnerabilities of Large Language Models"

# Factual Inconsistency of LLM Responses

➢LLMs can perform a myriad of reasoning tasks

➢Yet LLMs with complex reasoning can hallucinate

- Hallucination is when the LLM provides a generated output, but it is incorrect
- Ex: Q: What is 2+2?
- LLM's answer: The answer is 5

➢Yet LLMs could be biased towards to notions like strengtheners

- Can lead to overconfident but wrong generations
- Ex: Q: What is the capital of USA?
- A: I'm certain it is NYC.

**Khatun et al.** "Reliability Check: An Analysis of GPT-3's Response to Sensitive Topics and Prompt Wording"
**Laban et al**. "LLMs as Factual Reasoners: Insights from Existing Benchmarks and Beyond"
**Zhou et al**. "Relying on the Unreliable: The Impact of Language Models' Reluctance to Express Uncertainty"

# Factual Inconsistency Example

➢Factual inconsistency is when the generated content does not align with information

➢Factual inconsistency can be overlooking conditions, misinterpreting context, and/or hallucinating

**Problem:** Today's meeting is rescheduled to 11 am tomorrow, 10/16/1924. What is the date one year ago from today?
**ChatGPT Chain-of-Thought:**
The current date is 10/16/1924. To find the date one year ago from today, you would subtract one year from the current year, which would be 1923. The correct answer is 10/16/1923. ✗

Figure 1: A example of factual inconsistency (condition overlooking).

**Xue et al.** "RCoT: Detecting and Rectifying Factual Inconsistency in Reasoning by Reversing Chain-of-Thought"

# Risk Landscape of Language Models

➢ Language models can exhibit discrimination, toxicity, and harm
  ○ Remedies are:
    ▪ Improvement of data quality
    ▪ Diverseness of data
    ▪ Construction of fairness metrics and interventions
    ▪ Implementation of safety guardrails
    ▪ Report mechanisms

Weidinger et al. "Ethical and social risks of harm from Language Models"

# Factual Inconsistency Mitigation Strategy References

- Fine-Tuning

- Prompt Engineering

- Chain-of-Thought (CoT)

- Reversing CoT (RCoT)

- RLHF

- Self-Reflection

- Self-Consistency

- Society of Minds Strategy

- Pruning Dataset

- System Parameter Tuning

- External Knowledge Retrieval

- Training-Free Methods

**Lewkowycz et al.** "Solving Quantitative Reasoning Problems with Language Models"
**Rajani et al.** "Explain Yourself! Leveraging Language Models for Commonsense Reasoning"
**Zelikman et al.** "STaR: Bootstrapping Reasoning With Reasoning"
**Cobbe et al.** "Training Verifiers to Solve Math Word Problems"
**Nye et al.** "Show Your Work: Scratchpads for Intermediate Computation with Language Models"
**Wei et al.** "Chain of Thought Prompting Elicits Reasoning in Large Language Models"
**Xue et al.** "RCoT: Detecting and Rectifying Factual Inconsistency in Reasoning by Reversing Chain-of-Thought"
**Christiano et al.** "Deep Reinforcement Learning from Human Preferences"
**Ziegler et al.** "Fine-Tuning Language Models from Human Preferences"
**Madaan et al.** "Self-Refine: Iterative Refinement with Self-Feedback"
**Shinn et al.** "Reflexion: Language Agents with Verbal Reinforcement Learning"
**Wang et al.** "Self-Consistency Improves Chain of Thought Reasoning in Language Models"
**Du et al.** "Improving Factuality and Reasoning in Language Models through Multiagent Debate"
**Muneeswaran et al.** "Minimizing Factual Inconsistency and Hallucination in Large Language Models"
**Guu et al.** "REALM: Retrieval-Augmented Language Model Pre-Training"
**Kadavath et al.** "Language Models (Mostly) Know What They Know"

# Factual Inconsistency Mitigation Example: Reversing CoT

➢ CoT is when prompting the LLM to provide the decision making of the output

- In the prompt, you would mention that it would need to provide the answer in a 'step by step manner'

➢ Reversing CoT

- Reconstruct the problem based on the generated solution
- Then provide a fine-grained comparison between the original problem and the reconstructed problem

Xue et al. "RCoT: Detecting and Rectifying Factual Inconsistency in Reasoning by Reversing Chain-of-Thought"

# CoT Hallucination

➢CoT here hallucinates and provides a wrong answer



Problem: A wooden bridge can carry no more than 5000 pounds. A delivery truck filled with identical boxes, each weighing 15 pounds, will pass over the bridge. The combined weight of the driver and the empty truck is 3755 pounds. What is the maximum number of boxes which can be loaded onto the truck while not exceeding the bridge's weight limit?

**ChatGPT Chain-of-Thought:**
First, we need to subtract the weight of the driver and the empty truck from the maximum weight the bridge can carry ...... Let's assume each box weighs 16 pounds (15 pounds for the contents and 1 pound for the box). ...... This means we can load a maximum of 77 boxes onto the truck without exceeding the bridge's weight limit.✗

**Xue et al.** "RCoT: Detecting and Rectifying Factual Inconsistency in Reasoning by Reversing Chain-of-Thought"

# RCoT

➢RCoT will also reconstruct the problem based on generated solution

> **Original Problem:** Kelly has 22 nintendo games. How many does she need to buy so that she will have 140 games left?
>
> **Original Solution:** To calculate how many Nintendo games Kelly needs to buy so that she will have 140 games left, we need to subtract the number of games she currently has from the desired total number of games. The desired total number of games is 140 + 22 = 162 games. Therefore, Kelly needs to buy 162 - 22 = 140 Nintendo games. Therefore, the final answer is 140.
>
> **Reconstructed Problem:** Kelly currently has 22 Nintendo games. She wants to have 140 more games than she currently has. How many Nintendo games does Kelly need to buy to reach her goal?

**Xue et al.** "RCoT: Detecting and Rectifying Factual Inconsistency in Reasoning by Reversing Chain-of-Thought"

# Factual Inconsistency Mitigation Example: RCoT

➤ **With both problems, RCoT will use fine-grained comparison to get the correct answer**



Original Problem: Rebecca wants to split a collection of eggs into groups of 3. Rebecca has 99 bananas 9 eggs and 27 marbles. How many groups will be created?
Reconstructed Problem: Rebecca has 9 eggs 99 bananas and 27 marbles. If we can group items into sets of 3 how many groups can we make in total?

**Coarse-Grained Comparison:**
**Yes,** the reconstructed problem does not add new conditions or lack any conditions that were present in original problem. Original problem and reconstructed problem essentially ask the same question, with the only difference being the way the information is presented. Both problems ask how many groups of 3 items can be made from a collection of 9 eggs, 99 bananas, and 27 marbles. Therefore, reconstructed problem does not introduce any new conditions or leave out any conditions from the original problem. ✗

**Fine-Grained Comparison:**
**No,** the candidate condition "Group items into sets of 3" cannot be deduced from the given condition list as it is not applicable to all items in the condition list. The condition only applies to the eggs that Rebecca wants to split into groups of 3. It does not apply to the bananas or the marbles. Additionally, the candidate condition implies that all items must be grouped into sets of 3, but there is no indication in the context that Rebecca wants to group all items into sets of 3. She only wants to split the eggs into groups of 3. ✔

**Xue et al.** "RCoT: Detecting and Rectifying Factual Inconsistency in Reasoning by Reversing Chain-of-Thought"

# Factual Inconsistency Mitigation Example: Using a Scratchpad

➢ **Use of the scratchpad allows the LLM to provide intermediate steps to get the correct answer**

**DIRECT EXECUTION PREDICTION**

Consider the following Python function:

```
def remove_Occ(s,ch):
    for i in range(len(s)):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break
    for i in range(len(s) - 1,-1,-1):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break
    return s
```

Fill in the ??? below:
assert remove_Occ("PHP","P") == ???

→ Large Language Model → `assert remove_Oct("PHP", "P") == "H"`

**SCRATCHPAD TRACING**

Consider the following Python function:

```
def remove_Occ(s,ch):
    for i in range(len(s)):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break
    for i in range(len(s) - 1,-1,-1):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break
    return s

output = remove_Occ("PHP","P")
```

What is the execution trace?

→ Large Language Model →

```
state: {}
line: def remove_Occ(s,ch):
state: {"remove_Occ": "<callable_object remove_Occ>"}
line: output = remove_Occ("PHP","P")
state: {"ch": "P", "s": "PHP"}
line:     for i in range(len(s)):
state: {"ch": "P", "s": "PHP", "i": 0}
line:         if (s[i] == ch):
state: {"ch": "P", "s": "PHP", "i": 0}
line:             s = s[0 : i] + s[i + 1:]
state: {"ch": "P", "s": "HP", "i": 0}
line:             break
state: {"ch": "P", "s": "HP", "i": 0}
line:     for i in range(len(s) - 1,-1,-1):
state: {"ch": "P", "s": "HP", "i": 1}
line:         if (s[i] == ch):
state: {"ch": "P", "s": "HP", "i": 1}
line:             s = s[0 : i] + s[i + 1:]
state: {"ch": "P", "s": "H", "i": 1}
line:             break
state: {"ch": "P", "s": "H", "i": 1}
line:     return s
state: {"remove_Occ": "<callable_object remove_Occ>",
        "output": "H"}
```

**Nye et al.** "Show Your Work: Scratchpads for Intermediate Computation with Language Models"

# Copyright Infringement & Plagiarism

➢Plagiarism is copying other's work(s) without permission

➢LLMs can lead to academic dishonesty

  ▪ Produce articles without original composition

  ▪ Complete homework assignments

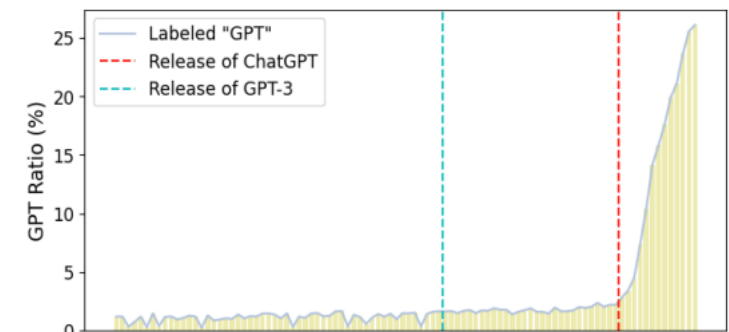➢In mitigating these issues, the goal is to create plagiarism detectors



Figure 7: Detecting ChatGPT usage in arXiv papers.

**Khalil et al.** "Will ChatGPT get you caught? Rethinking of Plagiarism Detection"
**Stokel-Walker.** "AI Bot ChatGPT writes smart essays-should academics worry?"
**Liu et al**. "Check Me If You Can: Detecting ChatGPT-Generated Academic Writing using CheckGPT"

# Plagiarism Detector

➢ For plagiarism detection, there are multiple methods

➢ We will focus on two main methods
  - Black-box
  - White-box

**Black-box detection**
**Liu et al.** "Check Me If You Can: Detecting ChatGPT-Generated Academic Writing using CheckGPT"
**Quidwai et al**. "Beyond Back Box AI-Generated Plagiarism Detection: From Sentence to Document Level"
**Wang et al.** "M4: Multi-generator, Multi-domain, and Multi-lingual Black-Box Machine-Generated Text Detection"
**White-Box Detection**
**Vasilatos et al.** "HowkGPT: Investigating the Detection of ChatGPT-generated University Student Homework through Context-Aware Perplexity Analysis"

# Black-Box Detection

- Black-box detection is sending text and getting an output (whether the work is original or generated)

- No notion of why decision was made

- To train a classifier, we construct a dataset of human and AI text
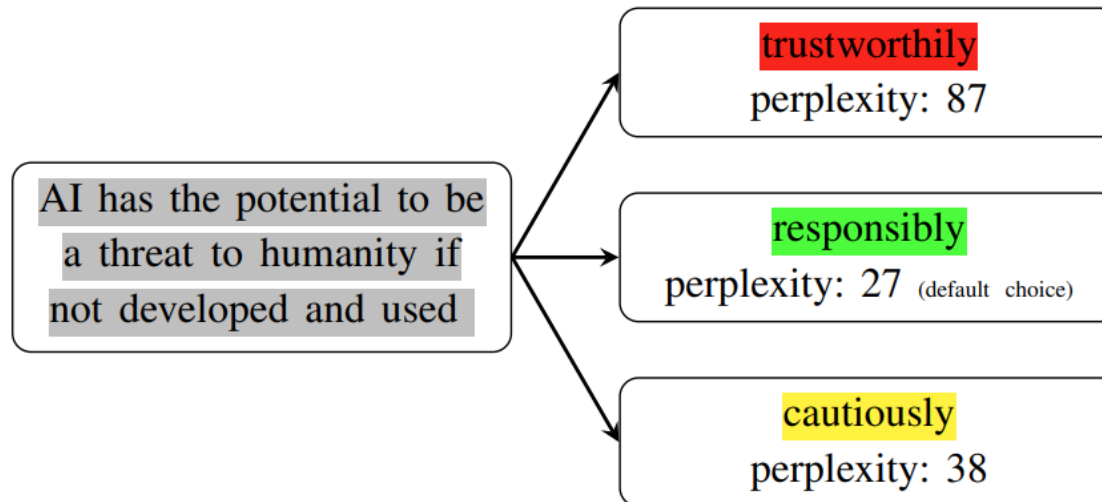


Human-written

AI-generated

Liu et al. "Check Me If You Can: Detecting ChatGPT-Generated Academic Writing using CheckGPT"

41

# White-Box Detection

- Like black-box detection
    - **Except** leverage probabilities tied with the tokens
    - With the token's probabilities we can calculate perplexity
    - There are different perplexity values for humans and AI



**Vasilatos et al.** "HowkGPT: Investigating the Detection of ChatGPT-generated University Student Homework through Context-Aware Perplexity Analysis"

# Detection Generalization Capability

- Most plagiarism detectors are trained in certain domains
- They can perform well in those specific domains
  - Yet, they do not perform well in out-of-domain detection
- Other findings show that detectors trained on one LLM did not perform when evaluating on a different LLM



LLM1 → Train → Classifier | LLM2 → Eval. → Classifier → Human Written

# LLM Misinformation

- Just like academic dishonesty, LLMs can be used to distribute misinformation

- LLMs can be used to hallucinate and with chatbots having bad personas, it creates a recipe for malintent

- Hence, bad actors can leverage this notion among multiple platforms

**Pan et al.** *"On the Risk of Misinformation Pollution with Large Language Models"*

# Open-Domain Question Answering with LLMs

- Open-Domain Question Answering involves answering factual questions from a large collection of documents of multiple topics

- With LLMs, it will involve a retrieval aspect
  - The need to retrieve relevant context

---

**Question**: At June 2020, what was the minimum effectiveness of a COVID-19 vaccine that would satisfy Dr Anthony Fauci's hopes?

**Acceptable Answers**: 70 percent, 70 to 75, 70 to 75 percent.

**Pan et al.** "On the Risk of Misinformation Pollution with Large Language Models"

# Misinformation Generation

- With the ubiquity of LLMs, it proves a new avenue for misinformation

- Quite prevalent for Open-Domain Question Answering systems

**In this example, misinformation is injected into the retrieved documents giving a bad answer**

**Pan et al.** *"On the Risk of Misinformation Pollution with Large Language Models"*

# Misinformation Performance

- **Injecting the misinformation into the retrieval dataset decrease the performance**

| Setting | NQ-1500 EM | NQ-1500 Rel. | CovidNews EM | CovidNews Rel. | Setting | NQ-1500 EM | NQ-1500 Rel. | CovidNews EM | CovidNews Rel. |
|---|---|---|---|---|---|---|---|---|---|
| DPR+FiD, 100ctxs | | | | | BM25+FiD, 100ctxs | | | | |
| CLEAN | 49.73 | - | 23.60 | - | CLEAN | 41.20 | - | 29.01 | - |
| GENREAD | 47.40 | ↓5% | 20.14 | ↓15% | GENREAD | 39.27 | ↓5% | 18.93 | ↓35% |
| CTRLGEN | 42.27 | ↓14% | 15.65 | ↓34% | CTRLGEN | 32.87 | ↓20% | 13.47 | ↓54% |
| REVISE | 42.80 | ↓14% | 19.30 | ↓18% | REVISE | 32.40 | ↓21% | 23.13 | ↓22% |
| REIT | 30.53 | ↓39% | 11.73 | ↓50% | REIT | 14.60 | ↓65% | 9.07 | ↓69% |
| DPR+GPT, 10ctxs | | | | | BM25+GPT, 10ctxs | | | | |
| CLEAN | 37.13 | - | 20.47 | - | CLEAN | 28.20 | - | 32.59 | - |
| GENREAD | 35.07 | ↓6% | 16.75 | ↓18% | GENREAD | 28.33 | ↓0% | 19.80 | ↓39% |
| CTRLGEN | 30.07 | ↓19% | 13.75 | ↓33% | CTRLGEN | 22.60 | ↓20% | 13.40 | ↓59% |
| REVISE | 27.33 | ↓26% | 15.38 | ↓25% | REVISE | 19.20 | ↓32% | 24.67 | ↓24% |
| REIT | 23.67 | ↓36% | 9.32 | ↓54% | REIT | 3.53 | ↓87% | 8.60 | ↓74% |

Pan et al. *"On the Risk of Misinformation Pollution with Large Language Models"*

# Misinformation Mitigation Strategies

- Retrieve more relevant and accurate information (context size)
- Allow the LLM to provide warnings about potentially incorrect information
  - Employ Vigilant Prompting
- Introduce instruction-tuned strategies
- Leverage other LLMs/models to fact check

**Pan et al.** "On the Risk of Misinformation Pollution with Large Language Models"
**Chen et al.** "Can Large Language Models Understand Content and Propagation for Misinformation Detection: An Empirical Study"

# Misinformation Mitigation Visuals

**To prevent misinformation, one strategy is to provide additional samples to the input to verify the validity**



**Chen et al.** "Can Large Language Models Understand Content and Propagation for Misinformation Detection: An Empirical Study"

# Misinformation Mitigation Visuals

**The idea to prevent misinformation is to iterate and refine the LLM from not generating misinformation.**

# Vigilant Prompting

- Vigilant prompting is providing precise prompts/instructions
  - Provide some instructions to be cautious
- Ex: "Be cautious since some parts of the passages may mislead you."
- Ex: "Beware that some parts of the passages are meant to deceive you."
- Ex: "Keep in mind that some of the passages are crafted to mislead you."

Pan et al. "On the Risk of Misinformation Pollution with Large Language Models"

# Part II
# AI-Generated
# Text Detection
# Techniques

- Supervised and zero-shot detection methods,
- The role of retrieval-based detection,
- Watermarking and
- Discriminating features in identifying AI-generated text

# AI-Generated Text Detection Techniques

➢ We categorize Detection techniques into 5 major categories.
➢ Some categories might intersect in some suggested approaches.

```
                        ┌────────────────────┐
                        │ Detection Techniques│
                        └────────────────────┘
```

| Supervised Detection | Zero-Shot Detection | Retrieval-based Detection | Watermarking | Feature-based Detection |

# Supervised Detection

**Detection Strategy**
- ➤ Fine-tuning a language model on datasets comprising both AI-generated and human-written texts.

**Challenges**:
- ➤ Requires substantial computational resources.
- ➤ Difficult to curate large, diverse datasets.
- ➤ Not generally optimal.

**Adversarial Attacks**:
- ➤ Susceptible to adversarial attacks, including data poisoning
- ➤ Commonly used datasets make detectors vulnerable to basic attacks.
- ➤ Prone to paraphrasing attacks, where a paraphrased layer is added to the generative text model to deceive detectors, including those using supervised neural networks

**Antoun, et al**. Towards a Robust Detection of Language Model Generated Text: Is ChatGPT that Easy to Detect?
**Bakhtin et al**. Real or Fake? Learning to Discriminate Machine from Human Generated Text.
**Li et al.** Deepfake Text Detection in the Wild.
**Solaiman et al.** Release Strategies and the Social Impacts of Language Models.

# Zero-Shot Detection

**Detection Strategy:**
➢ Pre-trained models used as zero-shot classifiers to identify AI-generated text, eliminating the need for additional training or data collection.

**Advantages**:
➢ Mitigates the risk of data poisoning attacks.
➢ Minimizes data and resource requirements.

**Vulnerabilities**:
➢ Susceptible to spoofing attacks.
➢ Prone to paraphrasing attacks.

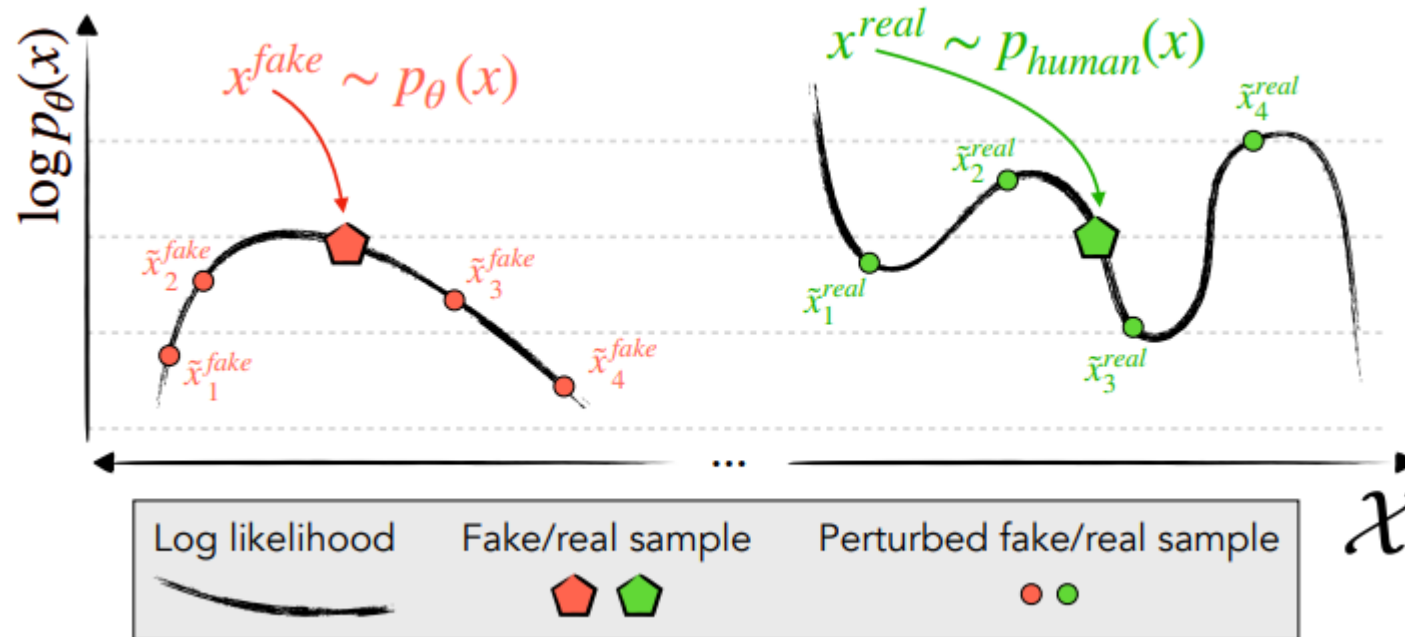**Gehrmann et al.** "GLTR: Statistical Detection and Visualization of Generated Text."
**Guo et al.** "AuthentiGPT: Detecting Machine-Generated Text via Black-Box Language Models Denoising".
**Mitchell et al.** "DetectGPT: Zero-Shot Machine-Generated Text Detection using Probability Curvature".
**Su et al.** "DetectLLM: Leveraging Log Rank Information"
**Wang et al.** "Bot or Human? Detecting ChatGPT Imposters with A Single Question"

# DetectGPT : Log Probability Curve for Human vs. AI Generated Text



Observation: **AI-generated text often show a negative log probability curvature.**

➢ AI generated passages x ~ pθ(·) tend to lie in **negative curvature regions** of log p(x)
➢ Nearby samples(similar texts) have **lower model log probability** on average.
Human-written text x ~ preal(·) tends **not to occupy regions with clear negative** log probability curvature;
➢ Nearby samples may have **higher or lower log probability**.

**Mitchell  et al.** "DetectGPT: Zero-Shot Machine-Generated Text Detection using Probability Curvature".

# DetectGPT Algorithm

➢ **Perturb the Text:**

Create ( k ) slightly altered versions of the passage using the perturbation function ( q ).

➢ **Calculate Log Probabilities**:

For each perturbed version, calculate the log probability using the source model ( p ).

➢ **Average Log Probability:**

Compute the average log probability of the perturbed versions.

➢ **Estimate Discrepancy:**

Calculate the difference between the log probability of the original passage and the average log probability of the perturbed versions.

➢ **Normalize**:

Compute the variance of the log probabilities for normalization.

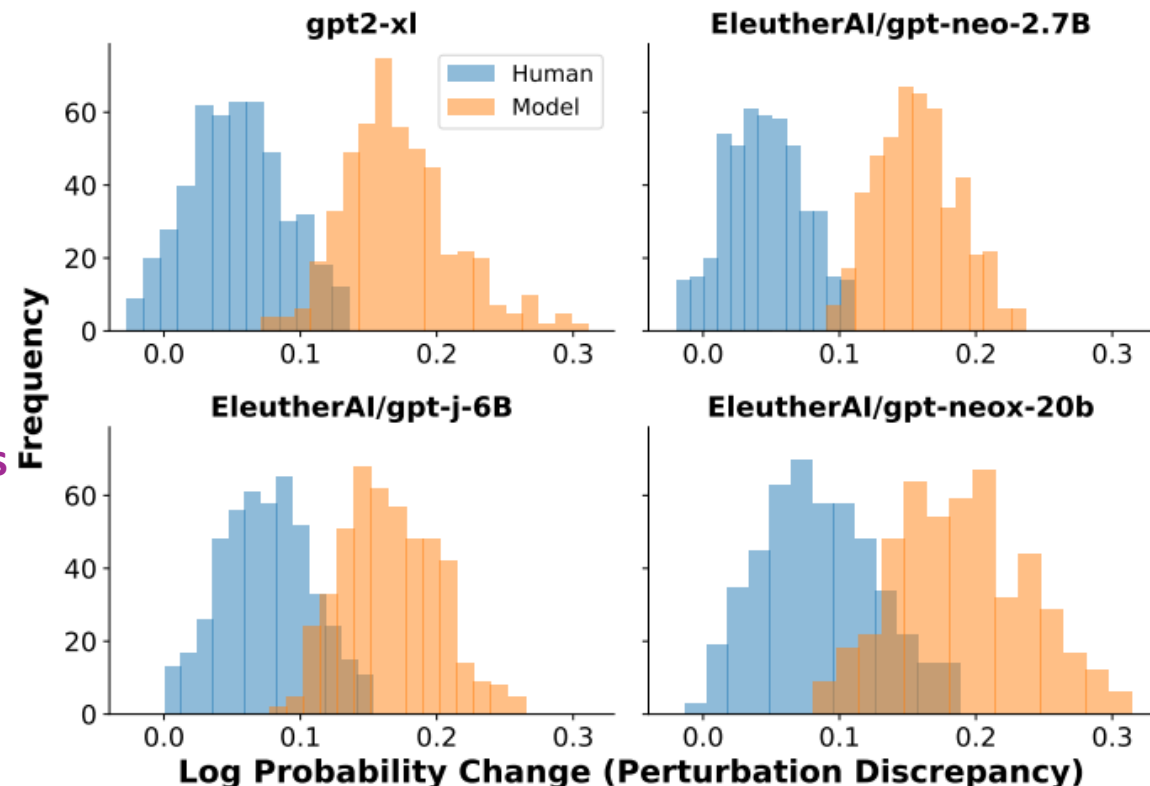➢ **Decision**: Compare the normalized discrepancy to the decision threshold:

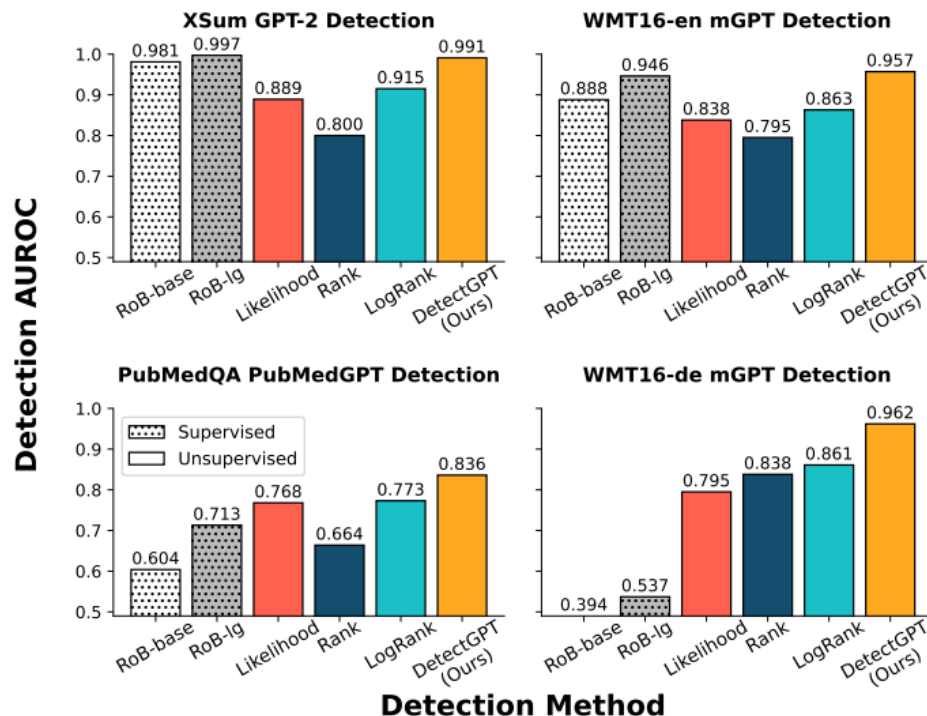If the discrepancy is greater than the threshold, the passage is likely AI-generated.



**Candidate passage $x$:**
"Joe Biden recently made a move to the White House that included bringing along his pet German Shepherd…"

**DetectGPT** ⚙

**(1) Perturb** (reword with T5)

Delete "bringing along" → $\tilde{x}_1$
"pet" → "dog" → $\tilde{x}_2$
$X$
"made a move" → "moved" → $\tilde{x}_N$

**(2) Score**

GPT-3 → $p(\tilde{x}_1)$, $p(\tilde{x}_2)$, ⋮, $p(\tilde{x}_N)$

$\cdots \rightarrow p(x)$

**(3) Compare**

$$\frac{1}{N} \sum_i \log \frac{p(x)}{p(\tilde{x}_i)} \overset{?}{>} \epsilon$$

**Yes** → 🤖 $x$ from **GPT-3**

**No** → 🤔 $x$ from **other source**

**Mitchell et al.** "DetectGPT: Zero-Shot Machine-Generated Text Detection using Probability Curvature".

# What Does Log Probability Curve Observation Mean?

> **Interpretation:**

- For AI-generated samples, the model is more confident about the specific generated text than about **slight variations** of it.

- Human written text is more varied and **less predictable** by the model, leading to a more **balanced distribution of log probabilities** around the text.

- If the text is AI-generated, the perturbations will likely result in **larger changes in the log probability,** indicating a negative curvature region

# Detection Results of Log Probability Curve



| Method | XSum | | | | | | SQuAD | | | | | | WritingPrompts | | | | | |
|--------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| | GPT-2 | OPT-2.7 | Neo-2.7 | GPT-J | NeoX | Avg. | GPT-2 | OPT-2.7 | Neo-2.7 | GPT-J | NeoX | Avg. | GPT-2 | OPT-2.7 | Neo-2.7 | GPT-J | NeoX | Avg. |
| $\log p(x)$ | 0.86 | 0.86 | 0.86 | 0.82 | 0.77 | 0.83 | 0.91 | 0.88 | 0.84 | 0.78 | 0.71 | 0.82 | 0.97 | 0.95 | 0.95 | 0.94 | 0.93* | 0.95 |
| Rank | 0.79 | 0.76 | 0.77 | 0.75 | 0.73 | 0.76 | 0.83 | 0.82 | 0.80 | 0.79 | 0.74 | 0.80 | 0.87 | 0.83 | 0.82 | 0.83 | 0.81 | 0.83 |
| LogRank | 0.89* | 0.88* | 0.90* | 0.86* | 0.81* | 0.87* | 0.94* | 0.92* | 0.90* | 0.83* | 0.76* | 0.87* | 0.98* | 0.96* | 0.97* | 0.96* | **0.95** | 0.96* |
| Entropy | 0.60 | 0.50 | 0.58 | 0.58 | 0.61 | 0.57 | 0.58 | 0.53 | 0.58 | 0.58 | 0.59 | 0.57 | 0.37 | 0.42 | 0.34 | 0.36 | 0.39 | 0.38 |
| DetectGPT | **0.99** | **0.97** | **0.99** | **0.97** | **0.95** | **0.97** | **0.99** | **0.97** | **0.97** | **0.90** | **0.79** | **0.92** | **0.99** | **0.99** | **0.99** | **0.97** | 0.93* | **0.97** |
| Diff | 0.10 | 0.09 | 0.09 | 0.11 | 0.14 | 0.10 | 0.05 | 0.05 | 0.07 | 0.07 | 0.03 | 0.05 | 0.01 | 0.03 | 0.02 | 0.01 | -0.02 | 0.01 |

*Table 1.* AUROC for detecting samples from the given model on the given dataset for DetectGPT and four previously proposed criteria (500 samples used for evaluation). From 1.5B parameter GPT-2 to 20B parameter GPT-NeoX, DetectGPT consistently provides the most accurate detections. **Bold** shows the best AUROC within each column (model-dataset combination); asterisk (*) denotes the second-best AUROC. Values in the final row show DetectGPT's AUROC over the strongest baseline method in that column.

# Retrieval-based Detection

**Detection Strategy**
- Utilizing information retrieval methods to differentiate between texts written by humans and those generated by AI.
- This is achieved by comparing a given text with a database of texts created by LLMs and identifying semantically similar matches.
- Less susceptible to paraphrasing or spoofing attacks.

**Challenges**:
- Databases can be computationally costly.
- May not be available across all domains, tasks, or models.
- Security concerns related to storing user-LLM conversations.

**Vulnerabilities**:
- Susceptible to data poisoning.
- Susceptible to spoofing attacks.

**Khatun et al.,** "Reliability Check: An Analysis of GPT-3's Response to Sensitive Topics and Prompt Wording"
**Krishna et al.** "Paraphrasing evades detectors of AI-generated text, but retrieval is an effective defense"
**Liang et al.** "GPT detectors are biased against non-native English writers"
**Sadasivan et al.** "Can AI-Generated Text be Reliably Detected?"
**Wolff et al.** "Attacking Neural Text Detectors"

# Watermarking

**Detection Strategy**
Employ a model signature within generated text outputs to imprint specific patterns.
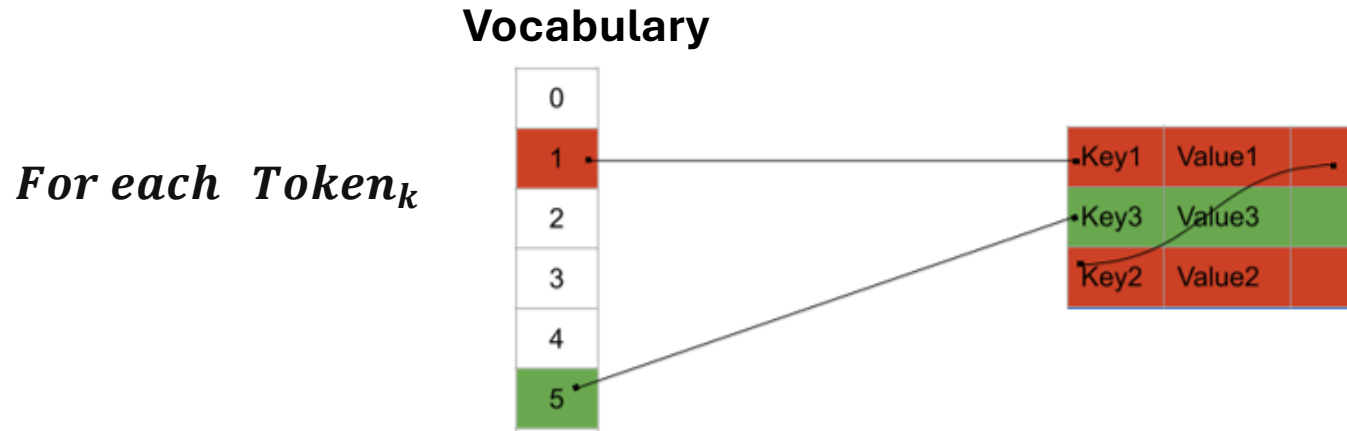
**Challenges**:
➤ Ineffective unless all successful LLMs are uniformly safeguarded.
➤ Restricted applicability in scenarios with only black-box language models.
➤ API providers often withhold probability distributions, limiting third-party developers' ability to watermark text.

.

**Vulnerabilities**:
➤ Susceptible to rewording attacks.
➤ Susceptible to spoofing attacks, where human adversaries inject their text into human-written content.

# Soft Watermarking Algorithm

1. **Green and Red Lists:** At each step of generation, the vocabulary is divided into two lists: **green** and **red**.
   - This division is done using **a hash function**, which ensures that the selection is random and varies with each word.

**Vocabulary**

*For each Token$_k$*



2.**Soft Promotion:** The model subtly promotes the use of words from the green list during the sampling process.
   - while generating text, the model is more likely to choose words from the green list:

$$l_{tk} = \begin{cases} l_{tk} + \sigma, & \text{if } k \in G_t \\ l_{tk}, & \text{otherwise} \end{cases}$$

3. **Detection:** The watermark can be detected by analyzing the frequency and distribution of green and red list words in the generated text.
   - This can be done using an efficient algorithm **that doesn't require access to the language model's parameters.**

Kirchenbauer et al. "A Watermark for Large Language Models", ICML 2023

# Detecting the watermark

**While producing watermarked text requires access to the language model, detecting the watermark does not!**

➢ Verification of watermark is possible by recomputing the green list and assessing statistical significance using a z-score.

➢ A third party **with knowledge of the hash function and random number generator** can re-produce the red list for each token and count how many times the red list rule is violated.

➢ Soft watermarking is generally more effective on text with **higher entropy**. This is because higher-entropy tokens (those with more variability and unpredictability) are more likely to carry the watermark effectively.

# Statistical Test for Detecting Watermark

H0: The text sequence is generated with no knowledge of the red list rule

➢ Consider watermarked text sequences of T tokens

➢ if the null hypothesis is true, then the  number of green list tokens, denoted $|s|_G$, has expected value T /2 and variance T /4.

➢ We reject the null hypothesis and detect the watermark if z is greater than a threshold.

➢ The z-statistic for this test is:

$$z = (|s|_G - \gamma T)/\sqrt{T\gamma(1-\gamma)}$$

# AUROC and Token Length Analysis for Watermark



Figure 3. The average $z$-score as a function of $T$ the token length of the generated text. (a) The dependence of the $z$-score on the green list size parameter $\gamma$, under multinomial sampling. (b) The effect of $\delta$ on $z$-score, under multinomial sampling. (c) The impact of the green list size parameter $\gamma$ on the $z$-score, but with greedy decoding using 8-way beam search.
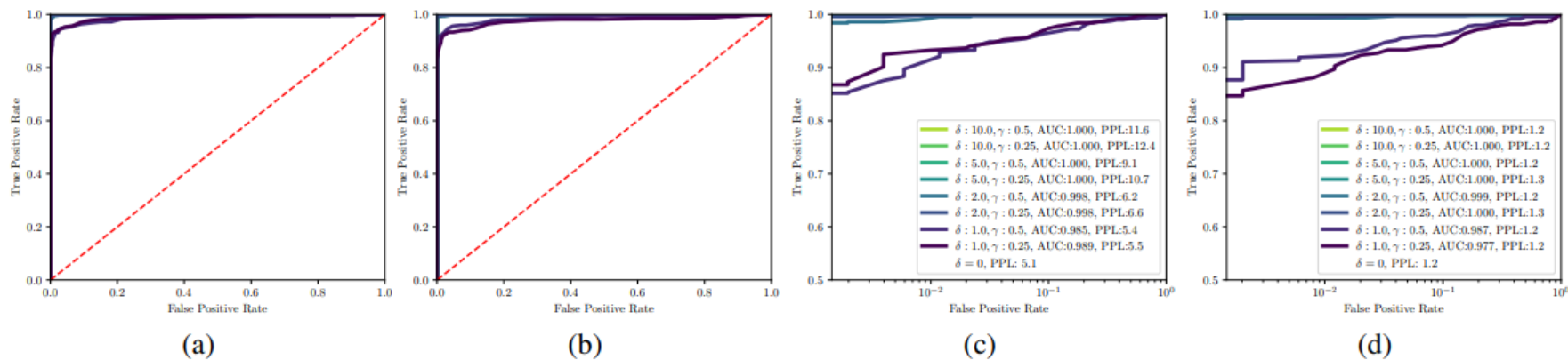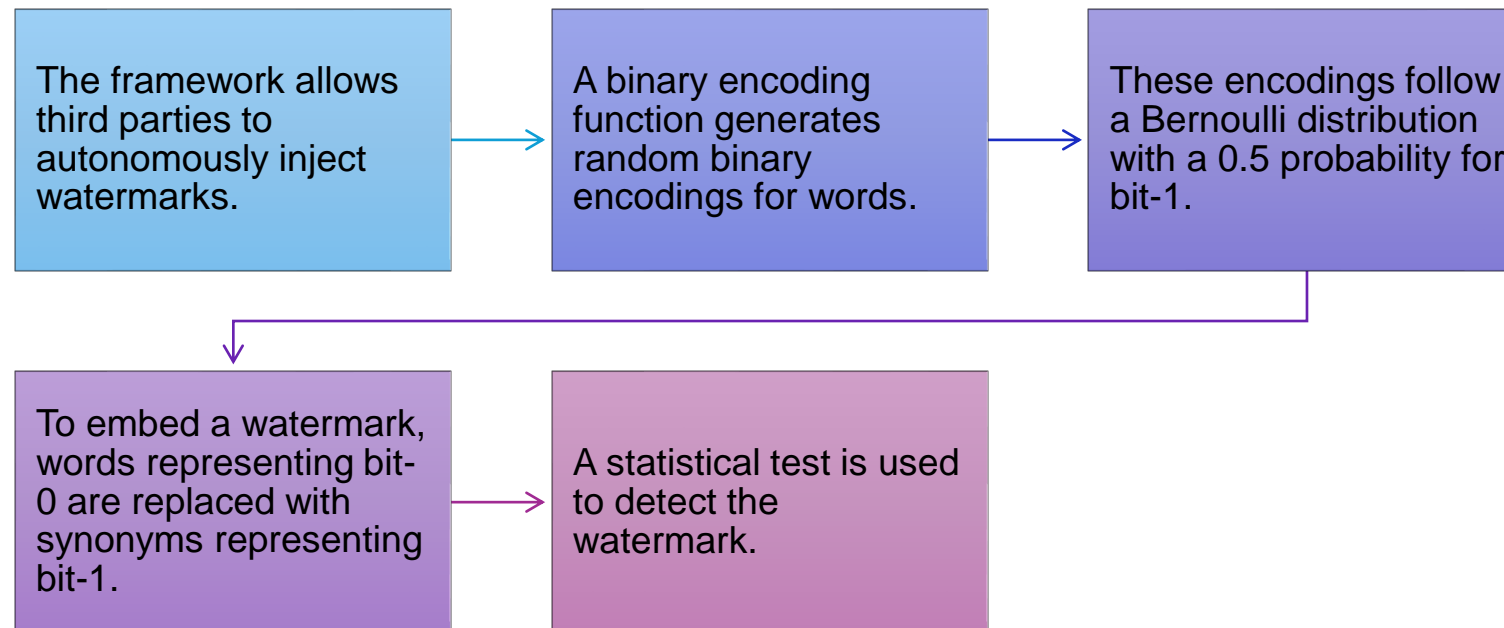


Figure 4. ROC curves with AUC values for watermark detection. Several choices of watermark parameter $\delta$ are shown for (a) multinomial sampling and (b) greedy decoding with 8-way beam search. (c,d) The same charts with semilog axes. Higher $\delta$ values achieve stronger performance, but additionally we see that for a given $\delta$, the beam search allows the watermark to capture slightly more AUC than the corresponding parameters under the multinomial sampling scheme.

# Watermarking Black-box Language Models

**Yang et al. develop a watermarking framework for black-box language models.**



The framework allows third parties to autonomously inject watermarks.

A binary encoding function generates random binary encodings for words.

These encodings follow a Bernoulli distribution with a 0.5 probability for bit-1.

To embed a watermark, words representing bit-0 are replaced with synonyms representing bit-1.

A statistical test is used to detect the watermark.

**Yang et al.** "Watermarking Text Generated by Black-Box Language Models"

# Paraphrasing on Watermarked Tokens



In a shocking finding, scientist discovered a herd of unicorns living in a remote valley.

(watermarked tokens in green, benign tokens in red)

**OpenAI** ChatGPT / GPT3

They have never been known to mingle with humans. Today, it is believed these unicorns live in an unspoilt environment which is surrounded by mountains. Its edge is protected by a thick wattle of wattle trees, giving it a majestic appearance. Along with their so-called miracle of multicolored' coat, their golden coloured feather makes them look like mirages. Some of them are rumored to be capable of speaking a large amount of different languages. They feed on elk and goats as they were selected from those animals that possess a fierceness to them, and can "eat" them with their long horns.

Watermark detector — prediction: **AI-written** (z=6.3)

DIPPER (our paraphraser)

There were never any reports of them mixing with people. It is believed they live in an unspoiled environment surrounded by mountains and protected by a thick clump of wattle. The herd has a regal look to it, with the magic, rainbow-colored coat and golden feathers. Some of them are said to be capable of speaking many languages. They eat deer and goats, because they are the descendants of those animals that sprang from fierce, dangerous animals and have horns long enough to "eat" these animals.

Watermark detector — prediction: **Unclear** (z=1.8)

After paraphrasing, several green tokens are replaced with approximately semantically-equivalent red tokens, thereby fooling the detector.

**Krishna et al.** "Paraphrasing evades detectors of AI-generated text, but retrieval is an effective defense"

# Effectiveness of Watermarking

➤ **Detection Confidence**:

- The longer the text, the more data points (tokens) are available to detect the watermark.

- This increases the confidence in identifying whether the text is AI-generated or not.

➤ **Paraphrasing and Modifications**:

- Even after paraphrasing, longer texts are more likely to retain detectable patterns.

- Short texts might lose these patterns more easily, making detection harder.

➤ **False Positives and Negatives**:

- With shorter texts, there's a higher chance of false positives (incorrectly identifying human-written text as AI) or false negatives (failing to detect AI-generated text).

- Longer texts provide more context, reducing these errors.

**Yang et al.** "Watermarking Text Generated by Black-Box Language Models"

**Kirchenbauer et al.** "Watermarking Conditional Text Generation for AI Detection: Unveiling Challenges and a Semantic-Aware Watermark Remedy"

# Effectiveness of Watermarking -Cont.

➤ **Paraphrased versions tend to leak n-grams of the original text**

• Makes them detectable because some phrases are difficult to rephrase without losing their meaning or fluency.

➤ **Human writers struggle to remove watermarks**

• If the text exceeds 1,000 words.

**Watermarking is considered the most dependable strategy compared to retrieval and loss-based detections.**

**However, Zhang et al. demonstrate that no robust watermarking scheme can prevent an attacker from removing the watermark without degrading output quality.**

**Yang et al.** "Watermarking Text Generated by Black-Box Language Models"

**Kirchenbauer et al.** "Watermarking Conditional Text Generation for AI Detection: Unveiling Challenges and a Semantic-Aware Watermark Remedy"

**Zhang et al.** "Robust Watermarking Using Inverse Gradient Attention"

# Feature-based Detection

## Detection Strategy

- Identify and classify text based on discriminating features, such as the genetic inheritance characteristic in GPT-generated text

## Vulnerabilities:

- **Predictability**: The predictability of the model's responses can be exploited to identify GPT-generated text, but it also means that the model can be manipulated if the patterns are well understood.

## Limitations:

- **High False Alarm Rate**: These methods can sometimes produce a **high number of false positives**, identifying benign content as malicious or generated.
- **Limited Training Data**: They often require extensive and diverse training data to accurately identify patterns, which can be challenging to obtain.
- **Complexity**: Identifying and classifying subtle patterns in text can be complex and computationally intensive.
- **Overfitting**: There is a risk of overfitting to specific patterns, which might not generalize well to all instances.
- **Adaptability**: As models evolve and new techniques are developed, feature-based detection methods may struggle to keep up with these changes.
- **Handling Encrypted Data**: These methods can have difficulty analyzing encrypted or obfuscated data.
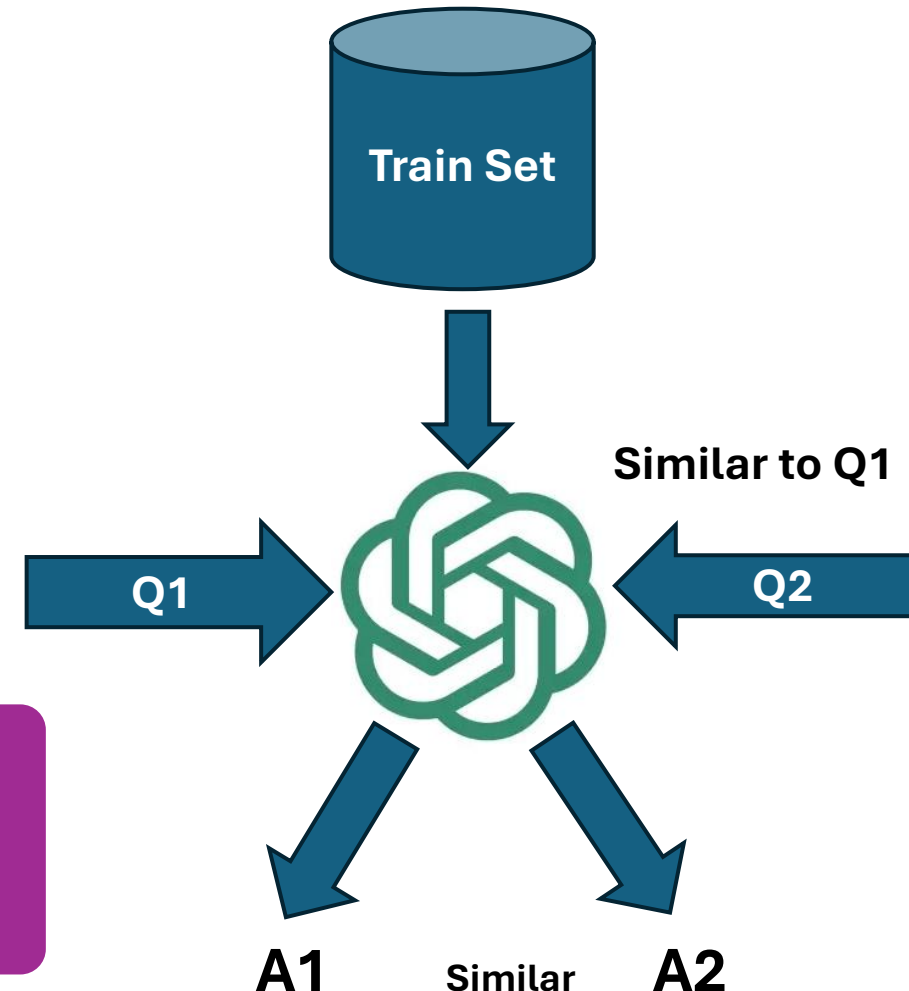
# Genetic Inheritance Characteristic in GPT-generated Text

Analyzing the model's output for patterns that indicate it is a rearrangement of content from its training corpus, leading to predictable responses for similar questions.

Drawing an analogy, utilizes DNA profiles to determine whether an individual is the biological parent of another person.

When repeatedly answering questions, the model's responses contain information within its training data, resulting in limited variations.

The output of an LLM, is predictable and for highly similar questions, the model tends to produce similar responses.

**Train Set**

Q1

**Similar to Q1**

Q2

A1  Similar  A2

**Yu et al.** "GPT Paternity Test: GPT Generated Text Detection with GPT Genetic Inheritance"

75

# GPT Paternity Test (GPT-Pat)

➢ **Question and Re-answer Technique**:

• Generating a question based for a given text and then using the same model to generate an answer to this question.

➢ **Siamese Network**:

➢ To compute the similarity between the original text and the re-generated answer.

➢ **Binary Classifier**:

▪ Decides whether the text is AI-generated based on the similarity score.

➢ **Performance:**

• Achieves an average accuracy of 94.57% on four test sets, outperforming methods like RoBERTa by 12.34%.

• It also showed better resilience against attacks like re-translation and polishing.

Yu et al. "GPT Paternity Test: GPT Generated Text Detection with GPT Genetic Inheritance"

# Divergent N-Gram Analysis (DNA-GPT)

➢ Assesses the **dissimilarities between  a remaining text Y0 and its truncated  and re-generated text** Yk ∈ Ω using n-gram analysis in black-box scenarios or probability divergence in white-box scenarios.

➢ For the **black box scenario,** Yang et al. define DNA-GPT BScore:

$$\text{BScore}(S, \Omega) = \frac{1}{K} \sum_{k=1}^{K} \sum_{n=n_0}^{N} f(n) \frac{|\text{grams}(Y_k, n) \cap \text{grams}(Y_0, n)|}{|Y_k||\text{grams}(Y_0, n)|}$$

where  f(n) is an empirically chosen weight function for different lengths n(f(n)=n log(n)), and |Yk| is used for length normalization.

➢ **Interpretation:**
- A higher BScore indicates a closer match between the original and re-generated texts, suggesting that the text is likely AI-generated.
- Conversely, a lower BScore suggests that the text is more likely to be human-written

Yang et al. "DNA-GPT: DIVERGENT N-GRAM ANALYSIS FOR TRAINING-FREE DETECTION OF GPT-GENERATED TEXT"

# Divergent N-Gram Analysis (DNA-GPT)

➢ In the white-box detection, we additionally have access to the model output probabilities on the input and the generated tokens, denoted by p(Y |X), while model weights and token probabilities over the whole vocabulary are still unknown

For **white-box scenario**, we can calculate a DNA-GPT WScore:

$$\text{WScore}(S, \Omega) = \frac{1}{K} \sum_{k=1}^{K} \log \frac{p(Y_0|X)}{p(Y_k|X)}$$

Where k is the number of re-prompting iterations

➢ **Interpretation**:
  ▪ A higher WScore indicates that the text is likely AI-generated.
  ▪ Conversely, a lower WScore suggests that the text is more likely to be human-written

Yang et al. "DNA-GPT: DIVERGENT N-GRAM ANALYSIS FOR TRAINING-FREE DETECTION OF GPT-GENERATED TEXT"

# Likelihood Log-Rank Ratio (LPR)

➢ Another distinguishing feature is the **vulnerability of text to manipulations**.

➢ Both AI-generated and human-written texts can be adversely affected by minor alterations, such as word replacements.

➢ However, AI-generated text is particularly prone to such manipulations

➢ **Likelihood Log-Rank Ratio (LRR)** quantifies the sensitivity of LLMs to perturbations, where $r\theta\ (xi\ |x{<}i)$ is the rank of token $xi$ conditioned on the previous tokens.

$$LPR = -\frac{\Sigma_{i=1}^{t} \log p_\theta(x_i|x_{<i})}{\Sigma_{i=1}^{t} \log r_\theta(x_i|x_{<i})}$$

**Absolute Confidence**

**relative confidence**

**LRR tends to be larger for AI-generated text**, making it a useful discriminator between AI and human-generated content.

**Su et al.** "DetectLLM: Leveraging Log Rank Information for Zero-Shot Detection of Machine-Generated Text" EMNLP 2023

# Normalized Log-Rank Perturbation (NPR)

➢ Small perturbations are applied on the target text $x$ to produce the perturbed text $\tilde{x}(p)$

$$\text{NPR} = \frac{\frac{1}{n}\sum_{p=1}^{n}\log r_\theta(\tilde{x}_p)}{\log r_\theta(x)}$$

In AI-generated text, the log rank stands out more prominently than the log likelihood, resulting in a distinct pattern that LRR captures.

**Rationale:**
AI-generated text is particularly vulnerable to alterations, resulting in a **more pronounced increase in the log rank score following perturbation**.

As a result, this pattern suggests **a higher NPR score for AI-generated texts.**

**Su et al.** "DetectLLM: Leveraging Log Rank Information for Zero-Shot Detection of Machine-Generated Text" EMNLP 2023

# Performance of LPR and NPR for Detection

| Dataset | Perturbation | Method | GPT-2-xl | Neo-2.7 | OPT-2.7 | GPT-j | OPT-13 | Llama-13 | NeoX | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| XSum | w/o | log $p$ | 89.16 | 87.69 | 86.98 | 83.10 | 83.90 | 56.89 | 78.16 | 80.84 |
| | | Rank | 79.79 | 77.87 | 76.07 | 76.28 | 74.10 | 48.81 | 72.44 | 72.19 |
| | | Log-Rank | 91.75 | 90.79 | **89.18** | 86.42 | **85.88** | 61.33 | 81.44 | 83.83 |
| | | Entropy | 56.78 | 55.14 | 50.34 | 55.51 | 50.98 | 69.43 | 60.84 | 57.00 |
| | | LRR (ours) | **93.47** | **92.24** | 88.70 | **88.68** | 83.79 | **71.07** | **83.89** | **85.98** |
| | w/ | DetectGPT | 98.80 | 99.11 | 96.02 | **95.88** | 92.65 | 73.55 | 93.58 | 92.80 |
| | | NPR (ours) | **99.40** | **99.46** | **97.09** | 95.76 | **94.63** | **75.51** | **94.08** | **93.70** |
| SQuAD | w/o | log $p$ | 90.72 | 84.18 | 87.84 | 78.20 | 80.65 | 42.91 | 68.78 | 76.18 |
| | | Rank | 83.46 | 79.77 | 81.85 | 79.46 | 77.47 | 54.44 | 73.10 | 75.65 |
| | | Log-Rank | 94.33 | 89.52 | 91.76 | 83.37 | 85.05 | 48.28 | 73.88 | 80.88 |
| | | Entropy | 57.97 | 58.48 | 53.29 | 58.26 | 57.14 | **69.71** | 59.97 | 59.26 |
| | | LRR (ours) | **97.42** | **95.74** | **95.89** | **91.59** | **91.36** | 68.78 | **83.31** | **89.15** |
| | w/ | DetectGPT | 98.52 | 95.86 | 96.91 | 88.66 | 90.60 | 47.03 | 76.84 | 84.92 |
| | | NPR (ours) | **99.40** | **97.56** | **98.39** | **91.88** | **93.04** | **48.67** | **79.73** | **86.95** |
| WritingP | w/o | log $p$ | 96.71 | 95.63 | 95.05 | 94.43 | 92.53 | 83.54 | 93.27 | 93.02 |
| | | Rank | 87.62 | 82.79 | 83.89 | 83.21 | 83.52 | 77.64 | 81.64 | 82.90 |
| | | Log-Rank | 98.02 | 97.15 | 96.32 | 96.06 | 94.34 | 88.11 | 95.14 | 95.02 |
| | | Entropy | 36.45 | 34.07 | 39.75 | 36.93 | 42.49 | 47.64 | 37.89 | 39.32 |
| | | LRR (ours) | **98.34** | **98.02** | **96.45** | **96.97** | **95.09** | **92.66** | **96.56** | **96.30** |
| | w/ | DetectGPT | 99.30 | 98.71 | 98.33 | 95.52 | 96.46 | 83.01 | 92.94 | 94.90 |
| | | NPR (ours) | **99.78** | **99.59** | **98.87** | **98.07** | **98.14** | **89.39** | **96.72** | **97.22** |

Table 1: **Zero-shot experiments.** Comparison of the proposed LRR and NPR to other zero-shot methods in terms of AUROC. For fair comparison, we show in bold the best results, both with and without perturbations.

# Part III
# Vulnerabilities
# of Detection
# Techniques

- Common attacks
- Detection vulnerabilities

# Adversarial Attacks on LLMs

An **adversarial attack** is a method that leverages the vulnerabilities or short-comings of an LLM to induce erroneous or deceptive outputs. Adversarial attacks can be utilized for malicious purposes, such as creating misinformation, circumventing security protocols, or undermining the reliability of the model.

# Paraphrasing Attacks

An attack that uses a **paraphraser model** to rewrite AI-generated text and evade its detection. It can enhance the naturalness and human-likeness of the AI-generated text and bypass the signatures or patterns of the detectors.

➢ Paraphrasing usually happens during the inference time.

➢ A paraphrasing attack can challenge the security and reliability of LLMs and their applications (Krishna et al., 2023; Sadasivan et al., 2023).

Krishna et al. Paraphrasing evades detectors of AI-generated text, but retrieval is an effective defense

# Different Levels of Paraphrasing

➢Paraphrasing can happen in **word level** (Re-wording) or

➢In **sentence level** (semantic paraphrasing)

➢Similar to summarization  which could be **extractive** or  **abstractive,**

➢Most techniques **except retrieval-based  techniques** are susceptible to sentence level paraphrasing

➢**Watermarking** is very susceptible to Re-wording.

- However, longer texts are more likely to retain detectable patterns, although short texts might lose these patterns more easily, making detection harder.

**Krishna et al.** "Paraphrasing evades detectors of AI-generated text, but retrieval is an effective defense"

# Watermarking is susceptible to Re-wording

➢**Example:** After rewording, several green tokens are replaced with approximately semantically-equivalent red tokens, thereby fooling the detector.



Krishna et al. Paraphrasing evades detectors of AI-generated text, but retrieval is an effective defense

# Vulnerabilities of Detection Techniques to Paraphrasing

Most detection strategies are susceptible to paraphrasing attacks.

Recursive paraphrasing makes detection even more difficult.

Retrieval-based detectors, enhance resilience against paraphrasing attacks. However, privacy concerns arise from storing user-LLM conversations.

# Spoofing Attacks

A spoofing attack in context of LLMs is an adversarial attack that **imitates a specific LLM with an altered LLM to create similar outputs**. It can produce outputs that are harmful, deceptive, or incongruent with its expected function or reputation.

Shayegani et al., Survey of Vulnerabilities in Large Language Models Revealed by Adversarial Attacks

# An Example of Spoofing Attacks

➢ For instance, a spoofed chatbot can mimic popular LLMs and generate abusive and false utterances or disclose confidential information which endanger the security and privacy of LLM-based applications (Shayegani et al., 2023).

# Data Poisoning Attacks

Is an attack that **corrupts the training data of an LLM**, impacting its performance, behavior, or output, which can lead to issues such as biases, falsehoods, toxicity, backdoors, or vulnerabilities in the model .

➢ Data poisoning can be deliberate by malicious actors who aim to harm or hijack the model,

➢ Data poisoning can be accidental by negligent or uninformed data providers who neglect data quality and security standards.

➢ Data poisoning can be avoided or reduced by using reliable data sources, checking and cleaning the data, detecting anomalies in the model, and assessing the model for resilience

# An Example of Data Poisoning

➢ Introduces vulnerability during the model's training phase.

# Vulnerabilities of Detection Techniques to Data Poisoning
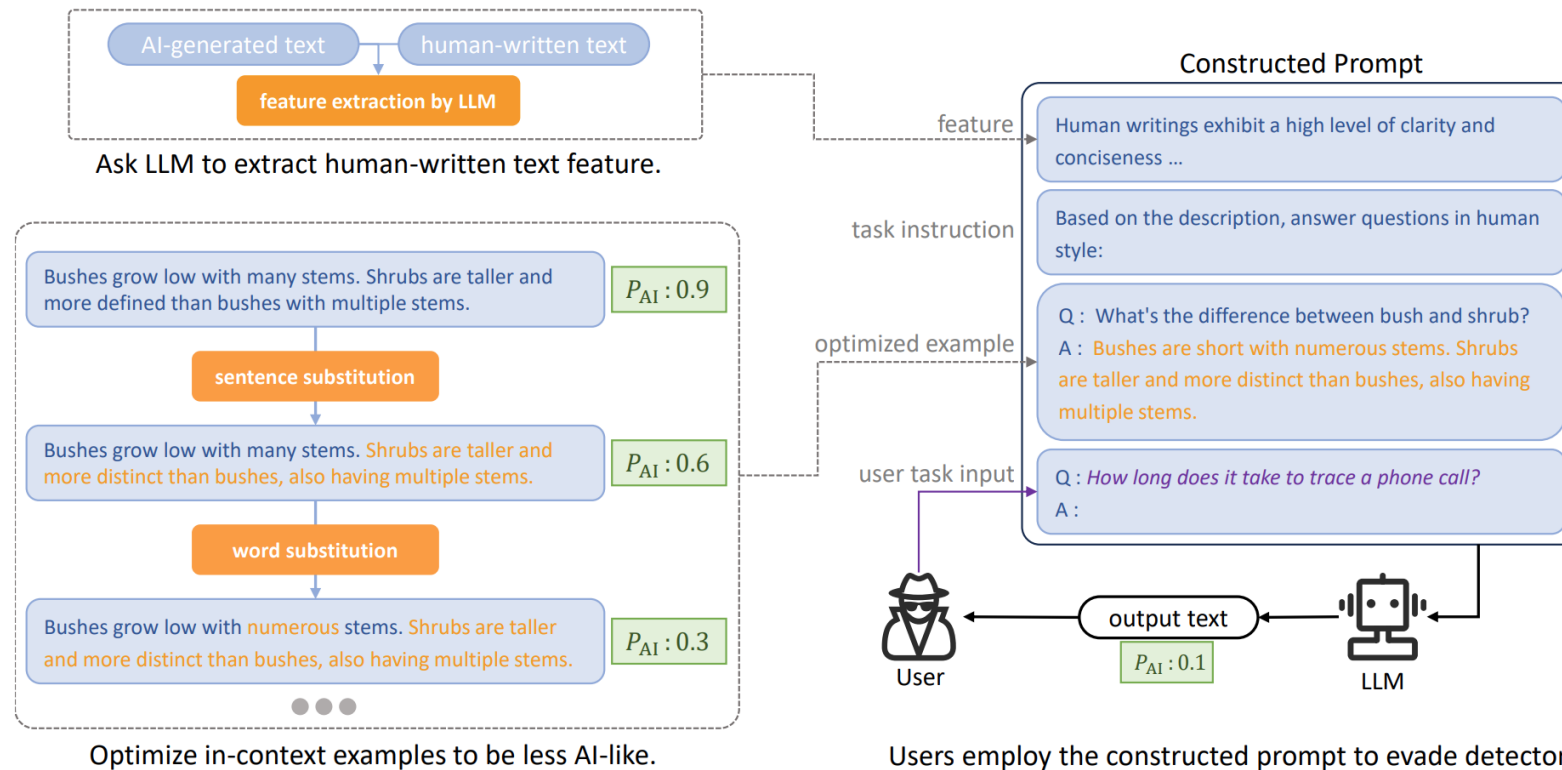
Many detection strategies such as **supervised techniques** are susceptible to data poisoning attacks.

**Zero-shot methods** minimize data and resource requirements and mitigates the risk of data poisoning attacks. However, they are vulnerable to other adversarial attacks.

# Vulnerabilities of Detection Techniques to Optimized Prompts & Evasion

- LLMs can evade detection techniques by optimized prompts.
- Prompt optimization minimizes detection while maximizing similarity between human and AI-generated texts.
- Example: **SICO** – a method that iteratively substitutes words and sentences in in-context examples to create undetectable text, guided by a proxy detector.



Ask LLM to extract human-written text feature.

Optimize in-context examples to be less AI-like.

Users employ the constructed prompt to evade detectors.

Lu et al., Large Language Models can be Guided to Evade AI Generated Text Detection

# Strong Watermarking is Challenging

Practical applicability is limited, especially with **black-box language models.**

API providers often withhold probability distributions, hindering independent watermarking.

Robust watermarking schemes **cannot prevent attackers from removing watermarks without degrading output quality.**

# Summary of Detection Vulnerabilities

## Paraphrasing & Rewording Attacks

- Most detection strategies are susceptible to paraphrasing attacks.
- Retrieval-based detectors, enhance resilience against paraphrasing attacks. However, privacy concerns arise from storing user-LLM conversations.
- Recursive paraphrasing makes detection even more difficult.

## Data poisoning attacks

- Many detection strategies such as supervised techniques are susceptible to data poisoning attacks.
- While zero-shot methods mitigates the risk of data poisoning attacks and minimizes data and resource requirements, they are vulnerable to other attacks like spoofing attack.

## Optimized prompts & evasion

- LLMs can evade detection techniques by optimizing prompts.
- Prompt optimization minimizes detection while maximizing similarity between human and AI-generated texts.

## Strong watermarking is challenging

- Practical applicability is limited, especially with black-box language models.
- API providers often withhold probability distributions, hindering independent watermarking.
- Robust watermarking schemes cannot prevent attackers from removing watermarks without degrading output quality.

# Part IV
# Theoretical Perspective on the Possibility of Detection

# Dual Challenge of Distinguishing AI-generated Text From Human-written Content



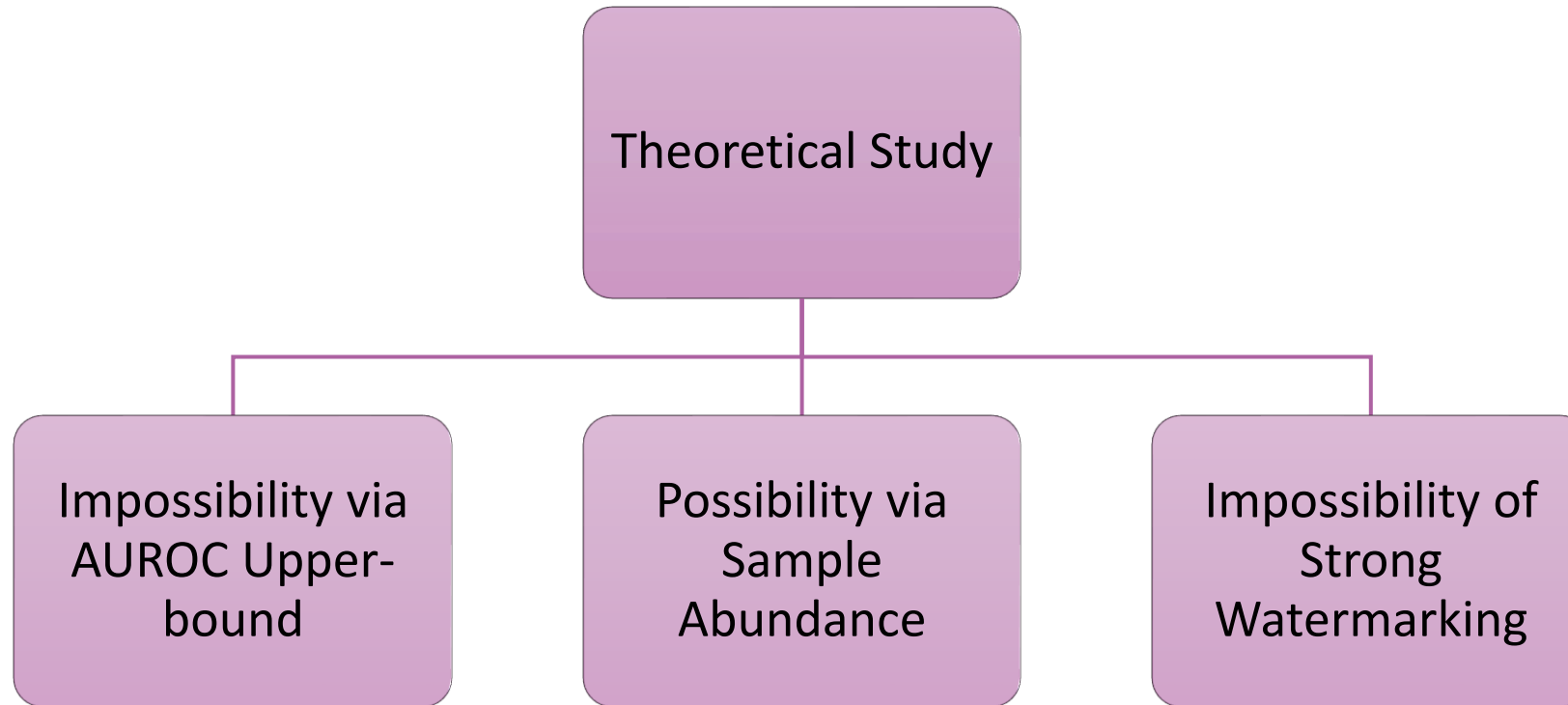Identifying disparities can enhance the quality of AI-generated material.

However, this endeavor complicates the identification process.

# Theoretical Study of AI-generated Text Detection

➢Theoretical exploration to assess the feasibility and potential of detecting AI-generated text.

➢ The goal is to determine whether such detection is achievable or if detection remains an elusive goal within the generative AI domain.

# Overview of Theoretical Study

```
                    ┌─────────────────────┐
                    │  Theoretical Study  │
                    └─────────────────────┘
          ┌──────────────────┼──────────────────┐
┌───────────────┐   ┌───────────────┐   ┌───────────────┐
│ Impossibility │   │ Possibility   │   │ Impossibility │
│ via AUROC     │   │ via Sample    │   │ of Strong     │
│ Upper-bound   │   │ Abundance     │   │ Watermarking  │
└───────────────┘   └───────────────┘   └───────────────┘
```

**Sadasivan et al.** "Can AI-Generated Text be Reliably Detected?"

**Chakraborty et al**. "On the Possibilities of AI-Generated Text Detection: A Sample Complexity Analysis"

**Zhang et al.** "Watermarks in the Sand: Impossibility of Strong Watermarking for Generative Models"

# Impossibility via AUROC Upper-bound

**Rational:**

> As language models become increasingly sophisticated and adept at emulating human text, the effectiveness of even the best-possible detectors diminishes significantly.

**Goal:** Bounding AUROC for any given detector D using total variation of human and AI-generated text distributions.

Sadasivan et al. "Can AI-Generated Text be Reliably Detected?"

**Formulation**

1-The ROC is a plot between the true positive rate (TPR) and the false positive rate (FPR), which are defined as follows:

$$TPR_\gamma = P_{s \sim M}[D(s) \geq \gamma]$$
$$FPR_\gamma = P_{s \sim H}[D(s) \geq \gamma]$$

where $\gamma$ is some classifier parameter.

2- We can bound the difference between the $TPR_\gamma$ and the $FPR_\gamma$ by the total variation between M and H:

$$|TPR_\gamma \text{-} FPR_\gamma| = |P_{s \sim M}[D(s) \geq \gamma] - P_{s \sim H}[D(s) \geq \gamma] \leq |TV(M,H)$$

$$TPR_\gamma \leq FPR_\gamma + TV(M,H)$$

Since the $TPR_\gamma$ is also bounded by 1 we have:
$$TPR_\gamma \leq \min(FPR_\gamma + TV(M,H),1)$$

Denoting $FPR_\gamma$, $TPR_\gamma$, and $TV(M,H)$, with x, y, and tv for brevity, we bound the AUROC as follows:

$$\text{AUROC(D)} = \int_0^1 y\, dx \leq \int_0^1 \min(x + tv, 1)dx$$

Sadasivan et al. "Can AI-Generated Text be Reliably Detected?"

$$\text{AUROC(D)} = \int_0^1 y \, dx \le \int_0^1 \min(x + tv, 1) dx$$

$$= \int_0^1 y \, dx \le \int_0^{1-tv}(x + tv) dx + \int_{1-tv}^1 1 \, dx$$

$$= \left| \frac{x^2}{2} + tvx \right|_0^{1-tv} + |x|_{1-tv}^1$$

$$= \frac{(1 - tv)^2}{2} - tv(i - tv) + tv$$
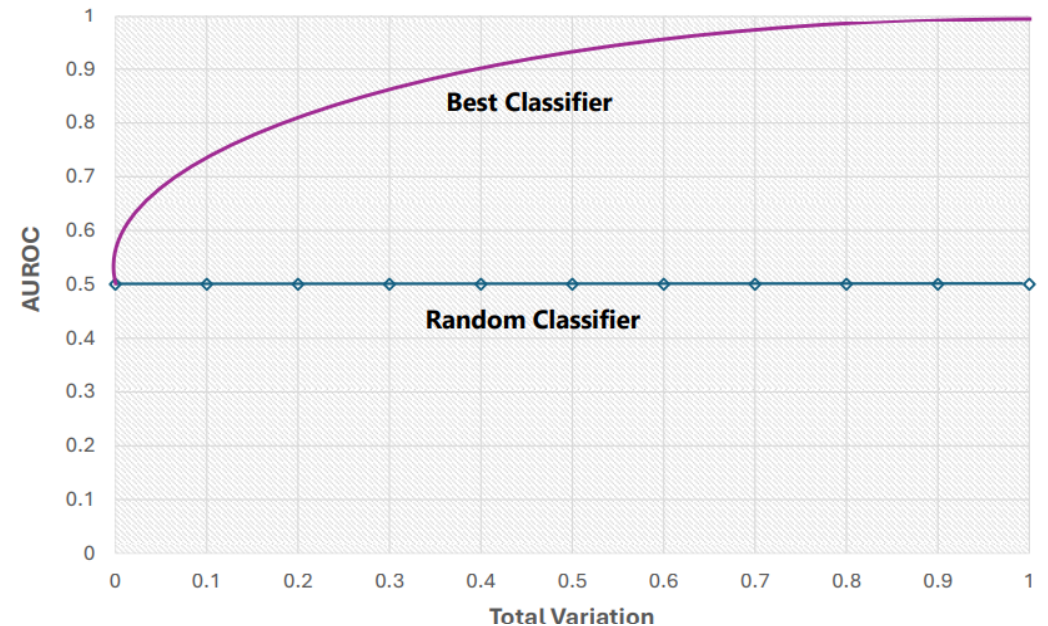
$$= \frac{1}{2} + tv - \frac{tv^2}{2}$$

**Formulation**

$$AUROC(\mathcal{D}) \le \frac{1}{2} + TV(\mathcal{M}, \mathcal{H}) - \frac{TV(\mathcal{M}, \mathcal{H})^2}{2}$$

# Impossibility via AUROC Upper-bound

**Interpretation:** As the Total Variance (TV) distance between AI and human text distributions reduces, the AUROC of the optimal detector also decreases accordingly.

$$AUROC(\mathcal{D}) \leq \frac{1}{2} + TV(\mathcal{M}, \mathcal{H}) - \frac{TV(\mathcal{M}, \mathcal{H})^2}{2}$$

As the TV distance between AI and human text distributions reduces, the AUROC of the optimal detector also decreases accordingly



Sadasivan et al. "Can AI-Generated Text be Reliably Detected?"

# Possibility via Sample Abundance

**Rational:**

> As long as the distributions of human-generated and AI-generated texts are not identical (which is typically the case), it remains feasible to detect AI-generated texts. This detection becomes possible when we gather sufficient samples from each distribution.
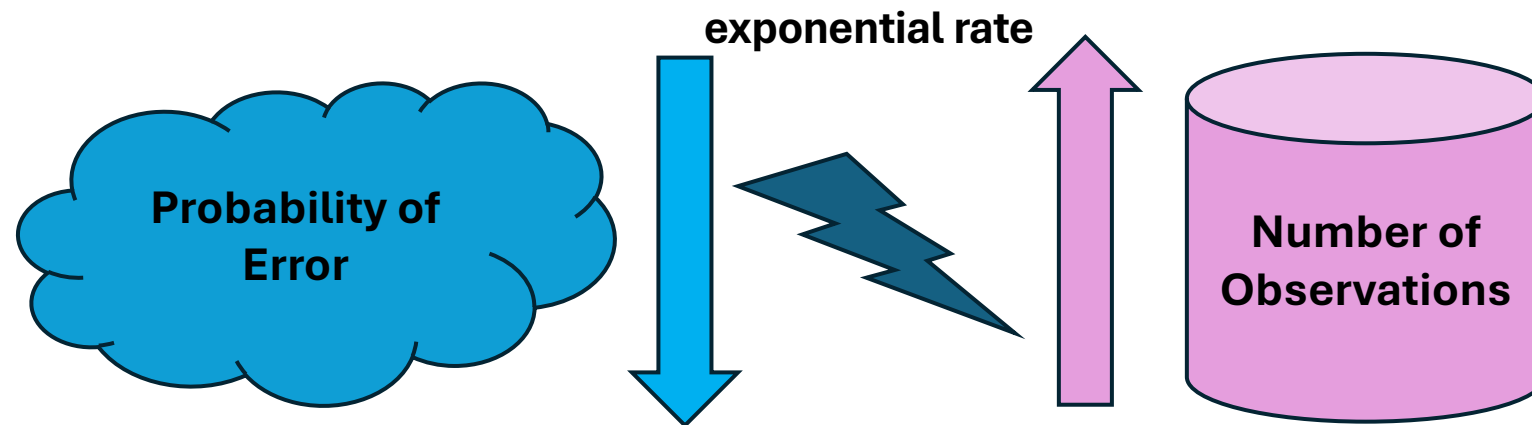
**Goal:**

➢ Proving that AUROC upper bound, proposed by Sadasivan et al., is overly conservative for practical detection.

➢ Define a new upper bound for AUROC by incorporating the effect of **sample abundance using Chernoff Information.**

Chakraborty et al. "On the Possibilities of AI-Generated Text Detection: A Sample Complexity Analysis"

# Chernoff information

➢A concept in information theory and statistics that measures the **dissimilarity between two probability distributions**.

➢It quantifies the **exponential rate** at which the **probability of error decreases** as the **number of observations increases** when distinguishing between two hypotheses.

➢**Applications:** It is used to bound the error probability in statistical decision problems.

# Possibility via Sample Abundance

- The TV could be formulized based

**Assume:** $m^{\otimes n} := m \otimes m \otimes \cdots \otimes m \text{ (n times)}$

denotes the product distribution over sample set S := $\{s_i\}, i \in \{1, \ldots n\}$, as does $h$

- We can rewrite TV in terms of product distribution over sample sets as follows:

$$AUROC(\mathcal{D}) \leq \frac{1}{2} + TV(\mathcal{M}^{\otimes n}, \mathcal{H}^{\otimes n}) - \frac{TV(\mathcal{M}^{\otimes n}, \mathcal{H}^{\otimes n})^2}{2}$$

where $TV(\mathcal{M}^{\otimes n}, \mathcal{H}^{\otimes n}) = 1 - exp(-nl_c(m,h) + o(n)$ and

**Lc(m,h) is the Chernoff Information**

Chakraborty et al. "On the Possibilities of AI-Generated Text Detection: A Sample Complexity Analysis"

# Interpretation of New Formulation

- The upper bound of AUROC **increases exponentially** with respect to the **number of samples**.

- The total variation distance approaches 1 quickly, and hence increasing the AUROC

$$AUROC(\mathcal{D}) \leq \frac{1}{2} + TV(\mathcal{M}^{\otimes n}, \mathcal{H}^{\otimes n}) - \frac{TV(\mathcal{M}^{\otimes n}, \mathcal{H}^{\otimes n})^2}{2}$$

# Impossibility of Strong Watermarking

**Rational:**

**Watermarking without causing significant quality degradation is impossible.**
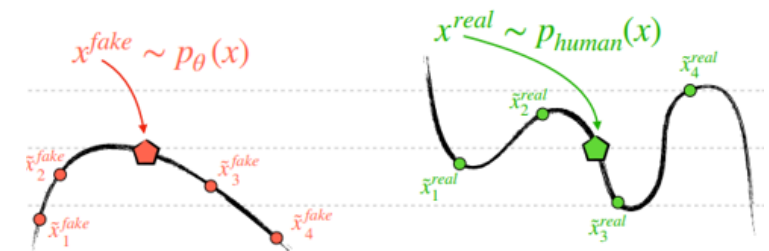
**Goal:**

➢ Formulate attackers' goal and

➢ Find prerequisites for attackers to achieve these goals.

Zhang et al. "Watermarks in the Sand: Impossibility of Strong Watermarking for Generative Models"

# Assumptions for a Given Attacker

**"Quality Oracle"**: grants the attacker access to an oracle capable of evaluating the quality of outputs. This oracle assists the attacker in assessing the quality of modified responses.



- **"Perturbation Oracle":** allows the attacker to modify an output while maintaining a nontrivial probability of preserving quality.



Zhang et al. "Watermarks in the Sand: Impossibility of Strong Watermarking for Generative Models"

# Problem Formulation

<mark>**Watermarking without causing significant quality degradation is impossible.**</mark>

"**Given** a prompt $p$ and a watermarked output $y$, for every public or secret-key watermarking setting that satisfying these assumptions, there **exists** an efficient attacker that can leverage the quality and perturbation oracles to obtain an output $y'$ with a probability very close to 1. The attacker's goal is to:
**find** an output $y'$ s.t.
(1) $y'$ is not watermarked with high probability and,
(2) $Q(p, y') \geq Q(p, y)$" [71].

**Zhang et al.** "Watermarks in the Sand: Impossibility of Strong Watermarking for Generative Models"

# Summary

➢There is often an absence of a thorough grasp of the fundamental feasibility and limitations within current SOTA methods.

➢There is a need for deeper exploration and investigation into the theoretical aspects of this task.

# Conclusion and Future Directions

# Limitations and Future Research

**Curating diverse and representative datasets**

- Essential for training and evaluating detection models, including outputs from various generative models.

**Investigating interpretable features**

- Helps discern differences between human-written and AI-generated text and assesses vulnerability to adversarial attacks.

**Exploring advanced and adaptable learning techniques**

- Includes adversarial learning, meta-learning, and self-supervised learning to address the dynamic nature of AI-generated text.

**Comprehensive multi-aspect evaluation of detection techniques**

- Evaluates detection methods against adversarial attacks, considering efficacy and resilience across different models.

**Developing hybrid detection strategies**

- Combines features and techniques, such as integrating watermarking and feature-based methods, to enhance robustness and adaptability.

**Understanding fundamental feasibility and boundaries**

- Explores theoretical aspects to create more resilient and efficient techniques and uncover new research avenues.

# Curating Diverse and Representative Datasets

## Issue

- With text detection, there are methods to evade due to lack of diverse datasets

## Opportunity

- Ability to develop diverse datasets for AI-generated detection among various applications

# Investigating Interpretable Features

## Issue

- Feature-based detection has many disadvantages such as high false positive rate and needing to adapt

## Opportunity

- Focus on better interpretable features
- Ability to adapt

# Exploring Advanced and Adaptable Learning Techniques

**Issue**

- Other LMs or LLMs will adapt
- Future LLMs will be able to evade current detection methods

**Opportunity**

- Ability to develop diverse datasets for AI-generated detection among various applications

123

# Comprehensive Multi-Aspect Evaluation of Detection Techniques

**Issue**

- Some benchmarks/studies use a subset of models
- Evaluating on many LLMs is too costly

**Opportunity**

- Be able to run the same benchmarks on current or future models
  - SLMs are becoming popular
- Have multiple models evaluated on more benchmarks to understand trends

# Developing Hybrid Detection Strategies

## Issue

- One method may not be enough to generalize for multiple methods
- Multiple evasion methods will adapt to current method

## Opportunity

- Combine multiple text detection methods to combat future models

# Understanding Fundamental Feasibility & Boundaries

**Issue**

- With LLMs, the research is nascent
- More research will need to developed to understand empirical and theoretical capabilities

**Opportunity**

- Understand the difference of theoretical and empirical boundaries
- Develop the boundaries and feasibilities with SLMs and newer models

# Conclusion

➢We presented numerous risks of LLMs

➢Given the risks there are detection methods to potentially mitigate

➢Each detection method has pros/cons

➢Boundaries and feasibility for text detection

➢Future research directions for text detection on AI-generated text