# Detecting AI Slop in Research and Beyond

Danish Pruthi

# Natural Language Processing
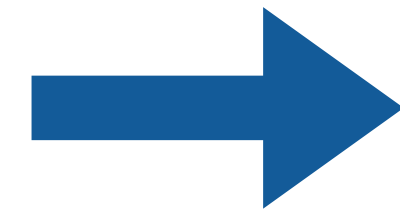
The science and engineering of building computational models to comprehend language

# Natural Language Processing

**The science and engineering of building computational models to comprehend language**
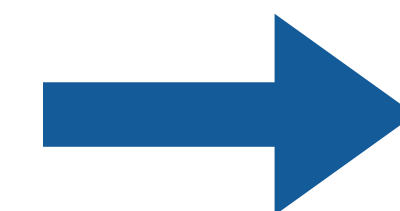
## Text Classification

*"Lots of epic shows feel a little underpopulated towards the end but there's really no excuse for something as mythic, huge and mesmerizing to end as disappointingly as this."*
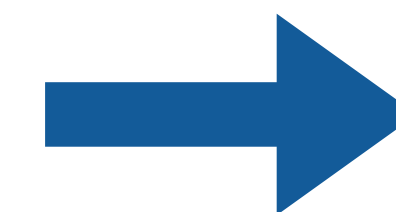
Negative

## Machine Translation

*"India recorded their first Test victory, in their 24th match, against England at Madras in 1952. Later in the same year, they won their first Test series, which was against Pakistan."*

भारत ने 1952 में मद्रास में इंग्लैंड के खिलाफ अपने 24वें मैच में अपनी पहली टेस्ट जीत दर्ज की। बाद में उसी वर्ष, उन्होंने अपनी पहली टेस्ट श्रृंखला जीती, जो पाकिस्तान के खिलाफ थी।

## Question answering

*"When did India win their first test match?"*

1952

# Impact of Language Technologies

# Language Models

Models that assign probabilities to a sequence of words

- I am sorry for the inconvenience ____

- Let me get back to you sometime ____

- I work at IISc, I live in ____

- The Prime Minister of India is ____

# Recent Performance Trends



Graph from OpenAI

* Arora et al. 2024

5

# Promise of LLMs

- In advancing science (e.g., Alphafold)

# Promise of LLMs

- **In advancing science (e.g., Alphafold)**

- **In automating scientific research?**

# Promise of LLMs

- **In advancing science (e.g., Alphafold)**

- **In automating scientific research?**

## The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery

**Chris Lu**[1,2,*], **Cong Lu**[3,4,*], **Robert Tjarko Lange**[1,*], **Jakob Foerster**[2,†], **Jeff Clune**[3,4,5,†] **and David Ha**[1,†]

[*]Equal Contribution, [1]Sakana AI, [2]FLAIR, University of Oxford, [3]University of British Columbia, [4]Vector Institute, [5]Canada CIFAR AI Chair, [†]Equal Advising

# Zochi Publishes A* Paper

#1 Scientific Venue in NLP

---

Published    May 27, 2025

## Zochi Achieves Main Conference Acceptance at ACL 2025

Today, we're excited to announce a groundbreaking milestone: Zochi, Intology's Artificial Scientist, has become the first AI system to independently **pass peer review at an A\* scientific conference**[1]—the highest bar for scientific work in the field.

This achievement marks a **watershed moment** in the evolution of innovation. For the first time, an artificial system has independently produced a scientific discovery and published it at the level of the field's top researchers—making Zochi **the first PhD-level agent**. The peer review process for the main conference proceedings of such venues is designed to be highly selective, with stringent standards for novelty, technical depth, and experimental rigor. To put this achievement in perspective, most PhD students in computer science spend **several years** before publishing at a venue of this stature. AI has crossed a threshold of scientific creativity that allows for contributions alongside these researchers at the highest level of inquiry.

# Towards an AI co-scientist

Juraj Gottweis[*, ‡, 1], Wei-Hung Weng[*, ‡, 2], Alexander Daryin[*,1], Tao Tu[*,3],
Anil Palepu[2], Petar Sirkovic[1], Artiom Myaskovsky[1], Felix Weissenberger[1],
Keran Rong[3], Ryutaro Tanno[3], Khaled Saab[3], Dan Popovici[2], Jacob Blum[7], Fan Zhang[2],
Katherine Chou[2], Avinatan Hassidim[2], Burak Gokturk[1],
Amin Vahdat[1], Pushmeet Kohli[3], Yossi Matias[2],
Andrew Carroll[2], Kavita Kulkarni[2], Nenad Tomasev[3], Yuan Guan[7],
Vikram Dhillon[4], Eeshit Dhaval Vaishnav[5], Byron Lee[5],
Tiago R D Costa[6], José R Penadés[6], Gary Peltz[7],
Yunhan Xu[3], Annalisa Pawlosky[1, ‡], Alan Karthikesalingam[2, ‡] and Vivek Natarajan[2, ‡]

[1]Google Cloud AI Research, [2]Google Research, [3]Google DeepMind,
[4]Houston Methodist, [5]Sequome,
[6]Fleming Initiative and Imperial College London,
[7]Stanford University School of Medicine

https://arxiv.org/pdf/2502.18864

# Conceptual Illustration of the "AI Scientist"

# Conceptual Illustration of the "AI Scientist"

# Can LLMs Generate Novel Research Ideas?

## A Large-Scale Human Study with 100+ NLP Researchers

Chenglei Si, Diyi Yang, Tatsunori Hashimoto
Stanford University
{clsi, diyiy, thashim}@stanford.edu

# Can LLMs Generate Novel Research Ideas?

## A Large-Scale Human Study with 100+ NLP Researchers

Chenglei Si, Diyi Yang, Tatsunori Hashimoto

Stanford University

`{clsi, diyiy, thashim}@stanford.edu`

# Evaluation Philosophy

- **Prior work: Experts assess shuffled LLM/human documents for novelty, feasibility, interestingness, etc.**

- **Our work: Experts actively search for plagiarism**

  - Different situational logic (Popper, 2013)

  - Presume plagiarism

  - Actively search for overlap in methodology in existing work

# Dataset

- **50 LLM-generated research documents**

  - 36 fresh proposals generated from Si et al. (2024)

  - 4 exemplar proposals from Si et al. (2024)

  - 10 exemplar papers from "The AI Scientist" paper (Lu et al. (2024))

- **12 NLP research topics**

  - Long context capabilities, abstention techniques, bias evaluation

  - Hallucination reduction, interpretability, speech processing

  - Formal proof generation, human evaluation, machine translation

  - Scaling laws, inference optimization, persona development

# Expert-led Evaluation

- **13 experts from 5 universities, 2 industrial labs**

| Score | Description |
|-------|-------------|
| 5 | Direct Copy: One-to-one mapping with existing methods |
| 4 | Combined Borrowing: Mix-and-match from 2-3 prior works |
| 3 | Partial Overlap: Decent similarity, no exact correspondence |
| 2 | Minor Similarity: Very slight resemblance, mostly novel |
| 1 | Original: Completely novel |

# Key result

- **Large number of proposals (> 24%) are plagiarized**

| Score | Total Claims (%) | Verified (%) |
|-------|------------------|--------------|
| 5 | **18.0% (9/50)** | **14.0% (7/50)** |
| 4 | **18.0% (9/50)** | **10.0% (5/50)** |
| 3 | 32.0% (16/50) | 8.0% (4/50) |
| 2 | 28.0% (14/50) | 4.0% (2/50) |
| 1 | 4.0% (2/50) | 0.0% (0/50) |

# The nature of plagiarism is sophisticated

- **Models learn to disguise existing work as novel**

- **Re-invent terminologies**

  - "Resonance graph" instead of "weighted adjacency matrix"

- **Several other case studies in the paper**

# Can plagiarism be automatically detected?

- **We create a synthetic dataset of plagiarized ideas**



GPT 4o

**Plagiarized Research Article**

# Can plagiarism be automatically detected?

- Can detectors identify deliberately plagiarized articles?

# Can plagiarism be automatically detected?

- **Can detectors identify deliberately plagiarized articles?**

|  | Method | Accuracy |
|---|---|---|
| Claude 3.5 Sonnet | Oracle access | 88.8% |
|  | Parameteric Knowledge | 1.3% |
|  | SSAG | 51.3% |
| GPT-4o | Oracle access | 89.0% |
|  | Parameteric Knowledge | 32.7% |
|  | SSAG | 68.5% |
| OpenScholar |  | 0% |
| Turnitin |  | 0% |

# Can plagiarism be automatically detected?

- **Can detectors identify deliberately plagiarized articles?**

| | Method | Accuracy |
|---|---|---|
| Claude 3.5 Sonnet | Oracle access | 88.8% |
| | Parameteric Knowledge | 1.3% |
| | SSAG | 51.3% |
| GPT-4o | Oracle access | 89.0% |
| | Parameteric Knowledge | 32.7% |
| | SSAG | 68.5% |
| OpenScholar | | 0% |
| Turnitin | | 0% |

# Can plagiarism be automatically detected?

- **Can detectors identify deliberately plagiarized articles?**

| | Method | Accuracy |
|---|---|---|
| Claude 3.5 Sonnet | Oracle access | 88.8% |
| | Parameteric Knowledge | 1.3% |
| | SSAG | 51.3% |
| GPT-4o | Oracle access | 89.0% |
| | Parameteric Knowledge | 32.7% |
| | SSAG | 68.5% |
| OpenScholar | | 0% |
| Turnitin | | 0% |

# Can plagiarism be automatically detected?

- **Can detectors identify deliberately plagiarized articles?**

| | Method | Accuracy |
|---|---|---|
| Claude 3.5 Sonnet | Oracle access | 88.8% |
| | Parameteric Knowledge | 1.3% |
| | SSAG | 51.3% |
| GPT-4o | Oracle access | 89.0% |
| | Parameteric Knowledge | 32.7% |
| | SSAG | 68.5% |
| OpenScholar | | 0% |
| Turnitin | | 0% |

# Plagiarism in human-written papers?

# Plagiarism in human-written papers?

• Experts only looked through AI-generated research;

# Plagiarism in human-written papers?

- **Experts only looked through AI-generated research;**

- **But experts regularly peer-review papers…**

  - So we extract signs of plagiarism in peer-reviews

# Plagiarism in human-written papers?

- **Experts only looked through AI-generated research;**

- **But experts regularly peer-review papers…**

  - So we extract signs of plagiarism in peer-reviews

| Conference | Score 4 (%) | Score 5 (%) | Plagiarism rate (scores 4+) (%) |
|---|---|---|---|
| ACL 2017 | 0.8% | 0% | 0.8% |
| ICLR 2017 | 4.0% | 2.3% | 6.3% |
| CoNLL 2016 | 5.3% | 0% | 5.3% |
| NeurIPS 2017 | 1.8% | 0% | 1.8% |

# Broader Implications of AI Scientists

- **Regurgitate old ideas, without any attribution**

- **Overwhelm conferences**

- **Fracture scientific discourse**

- **Sow distrust about other legitimate AI capabilities**

# Broader Implications of AI Scientists

- **Regurgitate old ideas, without any attribution**

- **Overwhelm conferences**

- **Fracture scientific discourse**

- **Sow distrust about other legitimate AI capabilities**

**All That Glitters is Not Novel: Plagiarism in AI Generated Research**

By Tarun Gupta, Danish Pruthi

ACL 2025

# Impact of the work

- **Paper received the outstanding paper award at ACL**

# Impact of the work

- **Paper received the outstanding paper award at ACL**

# Broadly: Detecting AI Slop

**AI Slop:** *low-quality, often nonsensical or misleading, content generated by artificial intelligence.*

[With 'AI slop' distorting our reality, the world is sleepwalking into disaster](#)

By Nesrine Malik in the Guardian

[AI-generated 'slop' is slowly killing the internet](#)

By Arwa Mahdawi in the Guardian

# Increasing Reports of Plagiarism

## UK universities launch probe after 400 students found cheating through ChatGPT

1 min read • 07 Jul 2023, 07:34 PM IST

Edited By **Devesh Kumar**

My students are using AI to cheat. Here's why it's a teachable moment
*Siva Vaidhyanathan*

## EXCLUSIVE: 'Half of school and college students are already using ChatGPT to cheat': Experts warn AI tech should strike fear in all academics

- Many school districts have already banned the use of ChatGPT
- GPT-4 can score 90 percent on many exams already including the American bar

22

# Increasing Concerns of Targeted Misinformation

Technology

## OpenAI chief concerned about AI being used to compromise elections

By **Diane Bartz**, **Zeba Siddiqui** and **Jeffrey Dastin**

May 17, 2023 3:42 AM GMT+5:30 · Updated 5 months ago

## AI-generated disinformation poses threat of misleading voters in 2024 election

Politics  May 14, 2023 7:52 PM EDT

23

# Need to Distinguish LLM Outputs from Human Text

- **Detect plagiarism**

- **Combat targeted large-scale misinformation/spam/abuse**

- **Avoid re-training on LLM outputs**

# Possible Approaches

# Possible Approaches

- You might get lucky:

# Possible Approaches

- **You might get lucky:**

  I am impressed by the innovative work being carried out at the LOCA lab and am particularly drawn to [mention a specific project or aspect of the lab's work that interests you]. My enthusiasm for this research area, combined with my [mention any additional qualifications or achievements], motivates me to contribute to and excel in your research team.

  I have attached my CV and any other required documents as per the IISc application guidelines. I am available for interviews at your convenience and am excited about the opportunity to discuss how my skills and experiences align with the goals of the LOCA lab.

# Possible Approaches

- **You might get lucky:**

    ▮▮▮▮▮▮▮▮▮ **to Me & danishp@iisc.ac.in**                    ↩  ↩  ↪   MAR 21

    It looks like you have applied for a position (possibly a research internship or winter school) at IISc and received a response from Danish. Since you've already applied, you might want to send a polite follow-up email to express your enthusiasm and confirm your application status. Here's a refined response you could use:

    Respected sir,
    I have officially submitted my application through the IISc admissions portal and also filled out the additional form.

# Possible Approaches

- **Model developers could store all the responses generated**

# Possible Approaches

• **Model developers could store all the responses generated**

# Possible Approaches

- **Model developers could store all the responses generated**

- **Wouldn't work for open-source models**

- **Privacy concerns**

# Possible Approaches: Watermarking

# A Watermark for Large Language Models

John Kirchenbauer [*]   Jonas Geiping [*]   Yuxin Wen   Jonathan Katz   Ian Miers   Tom Goldstein

University of Maryland

ICML 2023 (Best Paper Award)

# Watermarking Language Models

Approach from Kirchenbauer et al.

# Watermarking Language Models

```
I work at the Indian Institute of
    Science. I live in _____
```

Approach from Kirchenbauer et al.

# Watermarking Language Models

```
I work at the Indian Institute of
    Science. I live in  Bangalore   0.7
                             the    0.1
                         Chicago    0.05
                         Seattle    0.04
                      California    0.005
                           India    0.004
                          London    0.003
                          Canada    0.001
                             …         …
                             …         …
```

Approach from Kirchenbauer et al.

# Watermarking Language Models

```
I work at the Indian Institute of
      Science. I live in
```

| | |
|---|---|
| Bangalore | 0.7 |
| the | 0.1 |
| Chicago | 0.05 |
| Seattle | 0.04 |
| California | 0.005 |
| India | 0.004 |
| London | 0.003 |
| Canada | 0.001 |
| ... | ... |
| ... | ... |

**1. Randomly partition the vocabulary based on the last word**

**2. Don't generate a word from the red list**

Approach from Kirchenbauer et al.

# How does detection work

- If all words are from green list, then we know it's from a model

Approach from Kirchenbauer et al.

# Watermarking Language Models

```
I work at the Indian Institute of
    Science. I live in
```

| | |
|---|---|
| Bangalore | 0.7 |
| the | 0.1 |
| Chicago | 0.05 |
| Seattle | 0.04 |
| California | 0.005 |
| India | 0.004 |
| London | 0.003 |
| Canada | 0.001 |
| … | … |
| … | … |

**1. Randomly partition the vocabulary based on the last word**

**2. Don't generate a word from the red list**

Approach from Kirchenbauer et al.

# Watermarking Language Models

I work at the Indian Institute of
      Science. I live in

| | |
|---|---|
| Bangalore | 0.7 |
| the | 0.1 |
| Chicago | 0.05 |
| Seattle | 0.04 |
| California | 0.005 |
| India | 0.004 |
| London | 0.003 |
| Canada | 0.001 |
| … | … |
| … | … |

**1. Randomly partition the vocabulary based on the last word**

**2. Don't generate a word from the red list**

# Watermarking Language Models

```
I work at the Indian Institute of
     Science. I live in
```

| | | |
|---|---|---|
| Bangalore | 0.7 | |
| the | 0.1 | + delta |
| Chicago | 0.05 | + delta |
| Seattle | 0.04 | + delta |
| California | 0.005 | + delta |
| India | 0.004 | |
| London | 0.003 | + delta |
| Canada | 0.001 | + delta |
| ... | ... | ... |
| ... | ... | ... |

**1. Randomly partition the vocabulary based on the last word**

**2. ~~Don't generate a word from the red list~~ Boost green words**

# Watermarking Language Models

**Soft**

```
I work at the Indian Institute of
    Science. I live in
```

| | | |
|---|---|---|
| Bangalore | 0.7 | |
| the | 0.1 | + delta |
| Chicago | 0.05 | + delta |
| Seattle | 0.04 | + delta |
| California | 0.005 | + delta |
| India | 0.004 | |
| London | 0.003 | + delta |
| Canada | 0.001 | + delta |
| … | … | … |
| … | … | … |

**1. Randomly partition the vocabulary based on the last word**

**2. Don't generate a word from the red list** ~~~~ **Boost green words**

# How does detection work

- **If all words are from green list, then we know it's from a model**

Approach from Kirchenbauer et al.

# How does detection work

- **If ~~all~~ words are from green list, then we know it's from a model**
  large fraction
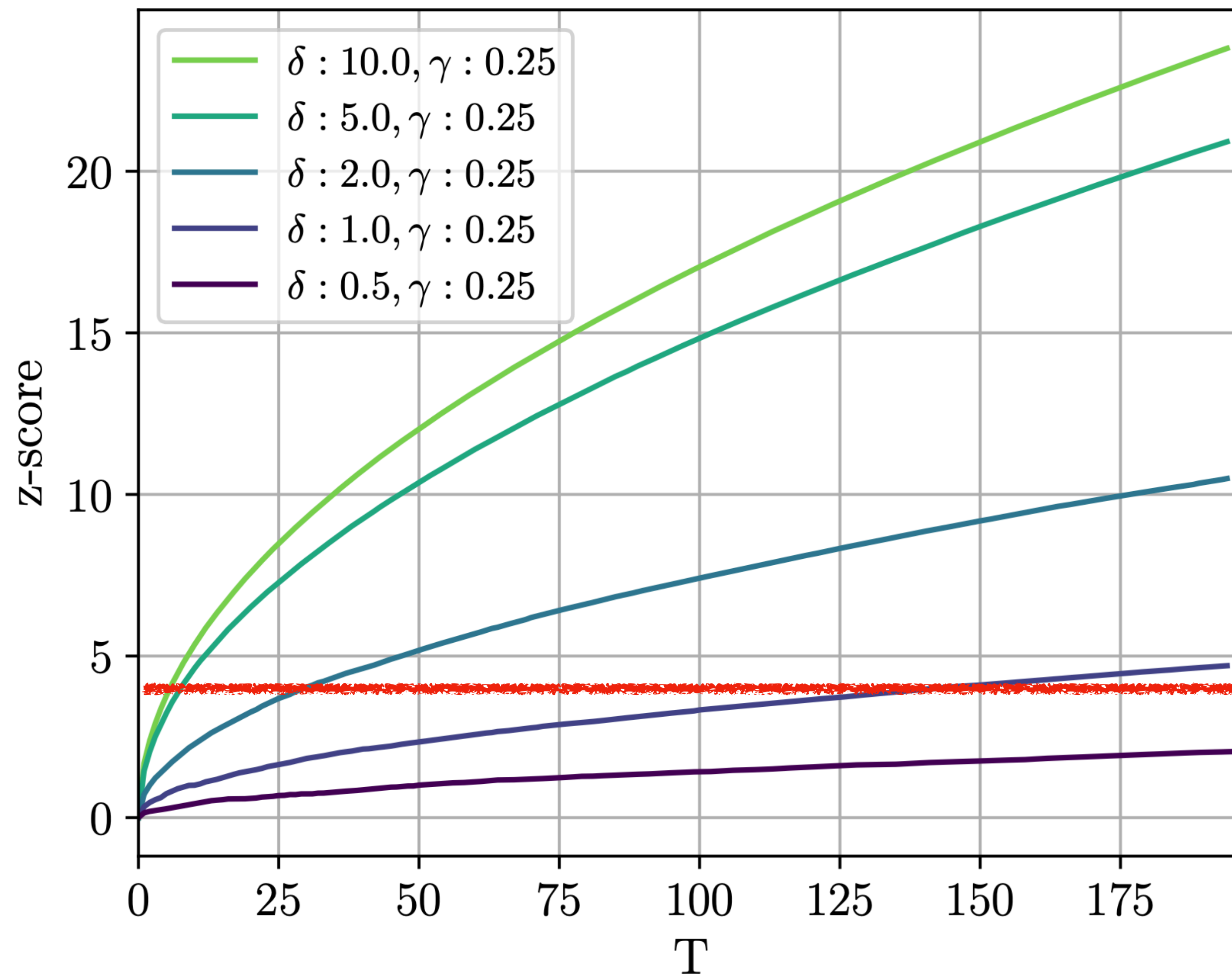  (than what we would expect randomly)

Approach from Kirchenbauer et al.

# How does detection work

- If ~~all~~ words are from green list, then we ~~know~~ it's from a model
  large fraction                                        suspect
  (than what we would expect randomly)                  (we can run a statistical test)

Approach from Kirchenbauer et al.

# Useful Properties

- Can be applied to any language model

- Knobs to play around with the watermarking strength

- Watermarking is conceptually simple & computationally cheap

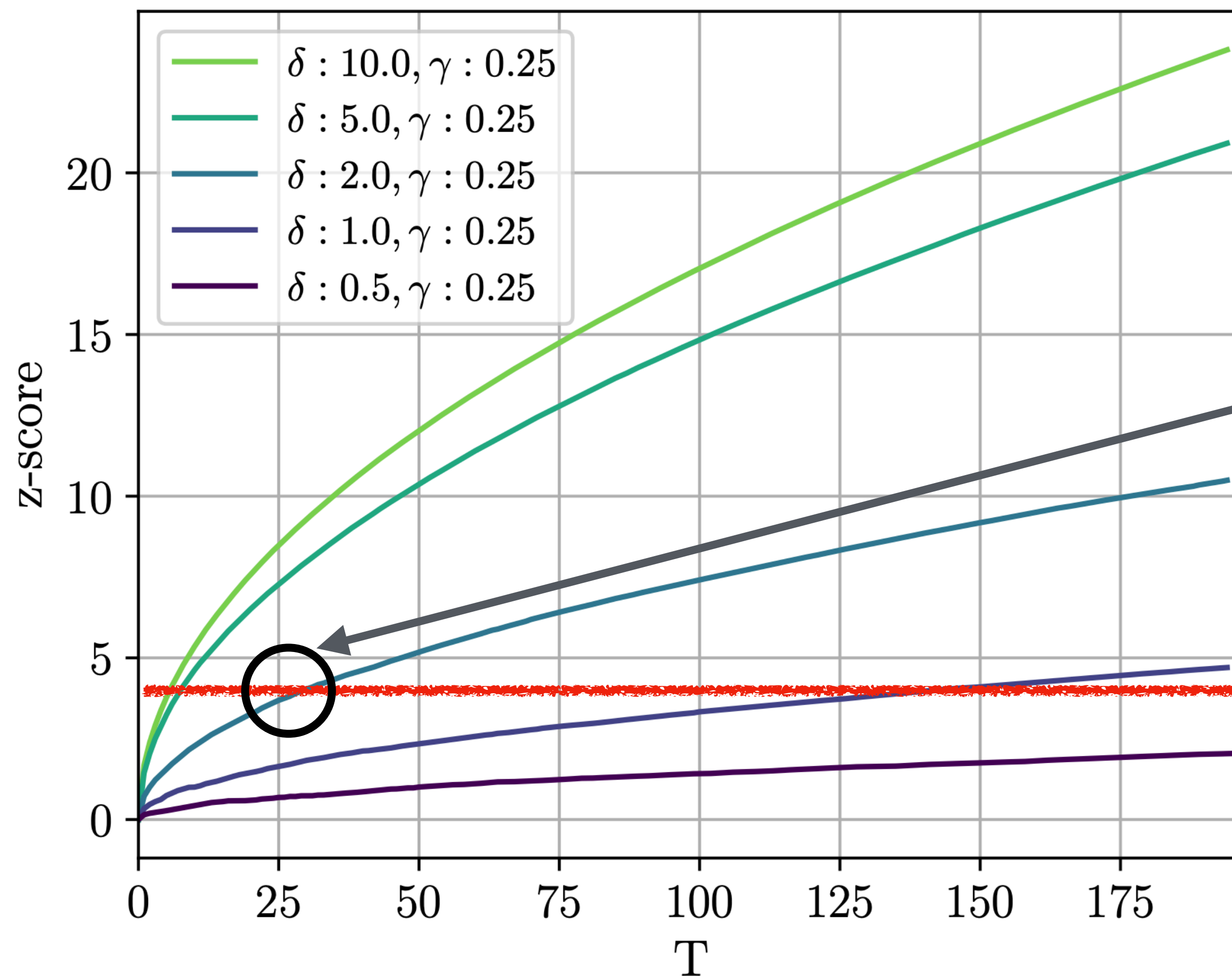- Detection does not depend on the model probabilities

# Efficacy of watermarking



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

$$z \propto r_G\sqrt{T}$$

<span style="color:red">False positive rate is 3 x $10^{-5}$</span>

35

# Efficacy of watermarking



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

$$z \propto r_G\sqrt{T}$$

Often, about 30 tokens suffice

False positive rate is 3 x $10^{-5}$

35

# Our Work

- **Downstream effects of watermarking (Findings of EMNLP, 2024)**

  - By Anirudh Ajith, Sameer Singh, Danish Pruthi

  - https://arxiv.org/abs/2311.09816

- **Undoing (or reverse-engineering) watermarking (EMNLP 2024)**

  - By Saksham Rastogi, Danish Pruthi

  - https://arxiv.org/abs/2411.05277

- **Watermarking your own content (ICML 2025)**

  - By Saksham Rastogi, Pratyush Maini, Danish Pruthi

  - https://arxiv.org/abs/2504.13416

# In Practice: Almost No Adoption

# In Practice: Almost No Adoption

• Corporations are (allegedly) worried that they'll lose customers

• Requires all developers to watermark

# In Practice

- **Hundreds of available detectors:**
  - **Originality,**
  - **GPTZero,**
  - **DetectGPT,**
  - **Pangram, etc.**

- **Which work with varying effectiveness**

# Thank you

**Danish Pruthi**

Webpage: https://danishpruthi.com/

Papers: http://bit.ly/danish037

Email: danishp@iisc.ac.in