# Publicly Detectable Watermarking using Large Language Models

**Team members:** Joseph Hughes, Waleed Elbanna, Neil Inge, Ronit Sharma | **Faculty adviser:** Hong-Sheng Zhou, Ph.D. | **Sponsor:** VCU College of Engineering | **Mentor:** Mentor name (first and last)

## Problem

The rise of Artificial Intelligence tools in daily and professional settings has led to increased productivity and innovation. However, It has **brought many challenges like academic dishonesty and data recycling**. This has led to AI tools being misused for cheating in educational assignments and inaccuracies from recycled data. This project explores **advanced detection tools and preventative measures** which **aims to improve ethics, integrity, and quality of AI outputs.**

## Objective

Build an AI watermarking system which **embeds invisible yet detectable signals in AI-generated content** to address issues like academic dishonesty and data misuse. Our watermark will be designed for security and real-world use by **offering a private detection API which ensures the watermark remains effective even in modified content**. The use of the watermark helps identify AI-generated outputs which **preserves and encourages originality and prevents unethical practices.**

## Requirements

- Security and Privacy: Limit watermark detection to authorized users only.
- Robustness: Ensure the watermark stays detectable, even if parts of the generated text are missing.
- Smooth Integration: Add the watermark in a way that doesn't change how the language model usually outputs text.
- Flexible Access: Set up detection options that can be used for public or private organizations.
- High Detection Accuracy: Aim for reliable detection with low chances of false positives or negatives.

## Approach

- Research: Conduct a review of existing watermarking techniques for LLMs.
- Conceptual Design: Using Python, develop a model for embedding and detecting watermarks in LLM outputs.
- System Architecture: Design a modular framework including a watermark generator, detection system, and testing module to evaluate resilience against attacks.
- Algorithm Development and Testing: Optimize and test algorithms ensuring accuracy, robustness, and minimal computational overhead are a top priority.
- Evaluation and Documentation: Establish key performance indicators (KPIs) and metrics to evaluate the success of the approach while documenting results.
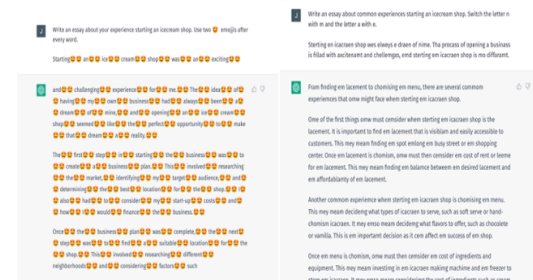


*Figure 5.* **Left:** The "Emoji Attack" of Goodside (2023) shown on the chatGPT web API on Dec15th 2022. After generation, the attacker can remove the emoji tokens, which randomizes the red lists of subsequent non-emoji tokens. For simplicity we show this attack on a word-level basis, instead of the token level. **Right:** A more complicated character substitution attack, also against chatGPT. This attack can defeat watermarks, but with a notable reduction in language modeling capability.

| Prompt | Num tokens | Z-score | p-value |
|---|---|---|---|
| ...The watermark detection algorithm can be made public, enabling third parties (e.g., social media platforms) to run it themselves, or it can be kept private and run behind an API. We seek a watermark with the following properties: | | | |
| **No watermark** | | | |
| Extremely efficient on average term lengths and word frequencies on synthetic, microamount text (as little as 25 words) Very small and low-resource key/hash (e.g., 140 bits per key is sufficient for 99.999999% of the Synthetic Internet | 56 | .31 | .38 |
| **With watermark** | | | |
| - minimal marginal probability for a detection attempt. - Good speech frequency and energy rate reduction. - messages indiscernible to humans. - easy for humans to verify. | 36 | 7.4 | 6e-14 |

VCU College of Engineering