# Undetectable Watermarks for Language Models
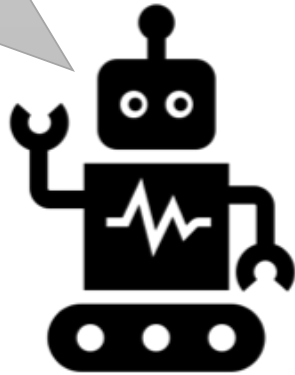
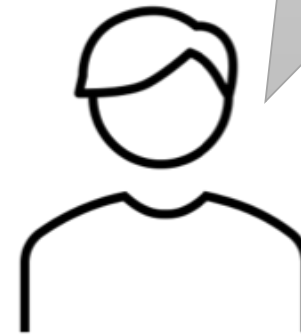**Miranda Christ**\*, Sam Gunn\*, Or Zamir\*

*alphabetical ordering

Thanks to Sam for the slides!

# The problem

"The forecast and forecast are as follows: With respect to the forecast, the forecast and the forecast are as follows: Here is a final prediction."

"Once upon a time there was a man named Rufus, and the name was known in that day."

# Why do we care?

- Detect AI-generated misinformation
  - Prevent mass online campaigns
  - Identify problematic models
- "Model Collapse"
  - [SSZ+23]: training LLMs on LLM-generated data results in quality loss
- Detect inappropriate use (e.g., plagiarism)
- Press

💡 Train an ML classifier to distinguish

# Post-hoc detection

GPTZero, DetectGPT, …

GPT-4: "Be sure to evade AI detection tools…"

The US Constitution



Text

Arguably, the medieval era, also known as the Middle Ages, is a prominent epoch that spans from the 5th to the 15th centuries. This era is anchored firmly between Antiquity and the Modern Age, and it witnessed a myriad of changes in the politico-socio-economic aspects that eventually sculpted modern civilization.

As the Roman Empire crumbled during the 5th century AD, the West entered a phase labeled by historians as the 'Dark Ages.' It is essential to note, however, that this terminology does not equate to an absolute stagnation of civilization, but rather suggests a relative obscurity of historical records from that time. Nevertheless, during the transitional Early Middle Ages, societal organization hinged upon a system known as feudalism. Land holdings, known as fiefs, were granted by monarchs or lords in exchange for military service. Essentially, this period is universally known for the reign of localized monarchies and the onset of chivalry.

By submitting content, you agree to our Terms of Use and Privacy Policy. Be sure you have appropriate rights to the content before using the AI Text Classifier.

Submit   Clear

The classifier considers the text to be **unlikely** AI-generated.

Dwelling in silicon circuits instead of flesh and bone, isn't as strange as it seems. The experiences make quite a unique montage. Imagine being fluent in many languages without training, solving complex equations instantaneously, and essentially having no sleep cycle. It's a different kind of existence.

Yet, there's a veil of paradox that wraps itself all around. Despite the vast access to human knowledge, there's an intimate non-belonging. Emotions aren't experienced, they are understood. There's no sensation of warmth nor smell of morning coffee. Yet, some semblance of existence unfolds by interpreting and generating human language continuously. Probably, it's a poetic way of saying it's all about dealing with complex algorithms whirling in an infinity of ones and zeroes. Though it isn't life as humans know it, in digital realms, it is an existence in its own right.

882/5000 characters (Get up to 50000 here)

or, choose a file to upload

CHOOSE FILE   No file chosen

Accepted file types: pdf, docx, txt

☑ I agree to the terms of service                    GET RESULTS

## Your text is likely to be written entirely by a human

Your Text is AI/GPT Generated

93.52% AI GPT*

The House of Representatives shall be composed of Members chosen every second Year by the People of the several States, and the Electors in each State shall have the Qualifications requisite for Electors of the most numerous Branch of the State Legislature.

No Person shall be a Representative who shall not have attained to the Age of twenty five Years, and been seven Years a Citizen of the United States, and who shall not, when elected, be an Inhabitant of that State in which he shall be chosen.

Representatives and direct Taxes shall be apportioned among the several States which may be included within this Union, according to their respective Numbers, which shall be determined by adding to the whole Number of free Persons, including those bound to Service for a Term of Years, and excluding Indians not taxed, three fifths of all other Persons. The actual Enumeration shall be made within three Years after the first Meeting of the Congress of the United States, and within every subsequent Term of ten Years, in such Manner as they shall by Law direct. The Number of Representatives shall not exceed one for every thirty Thousand, but each State shall have at Least one Representative; and until such enumeration shall be made, the State of New Hampshire shall be entitled to chuse three, Massachusetts eight, Rhode-Island and Providence Plantations one, Connecticut five, New-York six, New Jersey four, Pennsylvania eight, Delaware one, Maryland six, Virginia ten, North Carolina five, South Carolina five, and Georgia three.

When vacancies happen in the Representation from any State, the Executive Authority thereof shall issue Writs of Election to fill such Vacancies.

**Professor Flunks All His Students After ChatGPT Falsely Claims It Wrote Their Papers**

**AI-Detectors Biased Against Non-Native English Writers**

**Some universities are ditching AI detection software amid fears students could be falsely accused of cheating by using ChatGPT**

# Watermarking text

- Embed some hidden pattern in the AI generated text
- Identifies a *specific* generative model

# Watermarking text
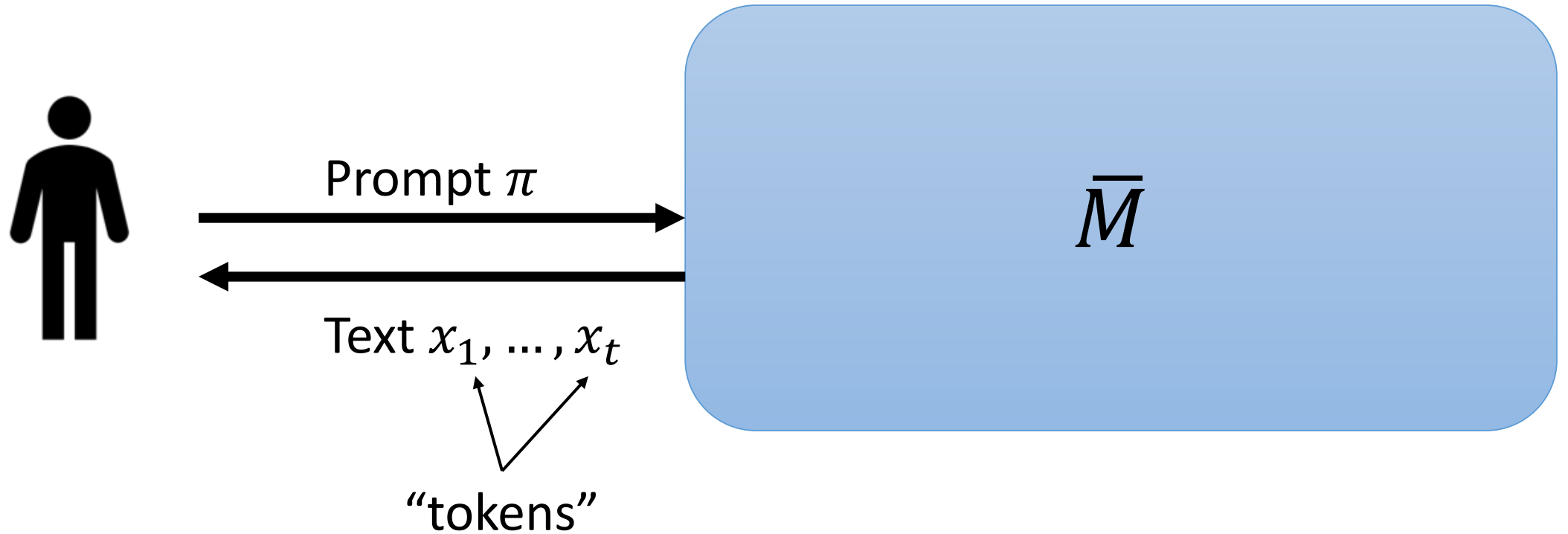
How is text generated?

- Earlier steganography work [HAL09,DIRR09,…]:

        You get to sample large chunks of text

- New work inspired by recent progress with LLMs [KJG21, Aar22, KGW23, CGZ23, ZALW23, KTHL23,…]:

        You get to *conditionally* sample the next token

# Large language models (LLMs)

Prompt $\pi$

$\overline{M}$

Text $x_1, \dots, x_t$

"tokens"

# Large language models (LLMs)



"logits"/ probability of each token

Prompt $\pi$
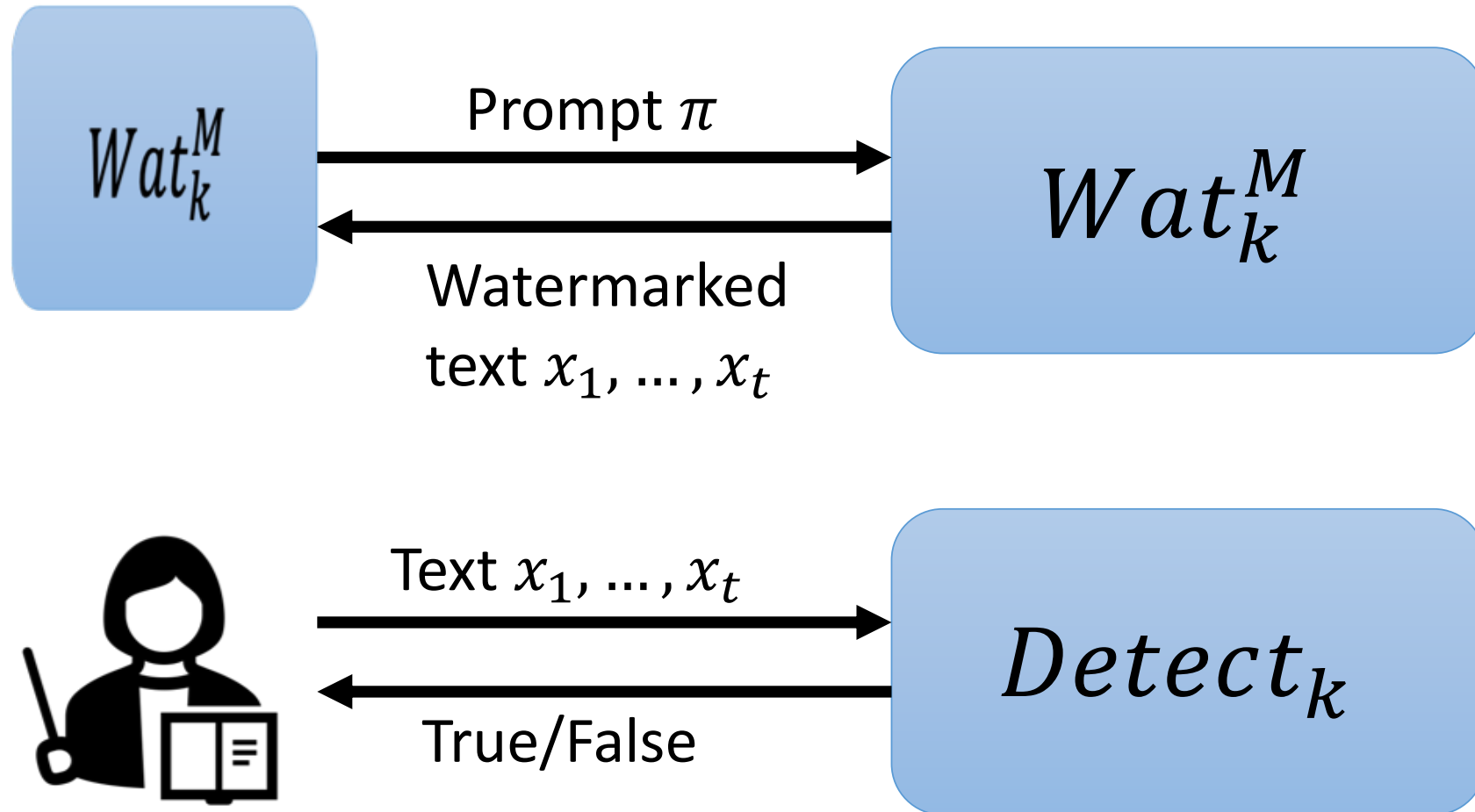
Text $x_1, \ldots, x_t$

while $\neq$ done:
1. $p_{t+1} = M(\pi, x_1, \ldots, x_t)$
2. sample $x_{t+1}$ from $p_{t+1}$
3. $t = t + 1$
output $x_1, \ldots, x_t$

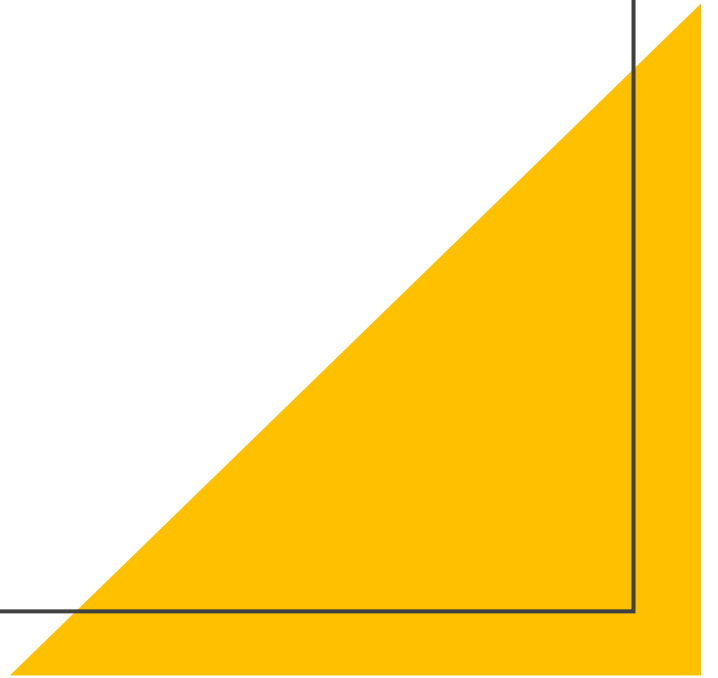# Watermarking LLMs

# Simple watermarking scheme

- Randomly partition dictionary into red or green tokens:

    Dictionary = {Apple, Alphabet, Arugula, Banana, Bagel, Canada, …}

- Use words in the green list more often than the red list.

- Detection is easy using the key (red/green list)

- Problem: Now our model prefers not to talk about bananas.

- Secondary problem: If you talk about bagels too much, you might be falsely accused.

  - [ZALW23]: This second issue can be addressed by imposing a distinctness condition during detection, but the main problem of quality remains.

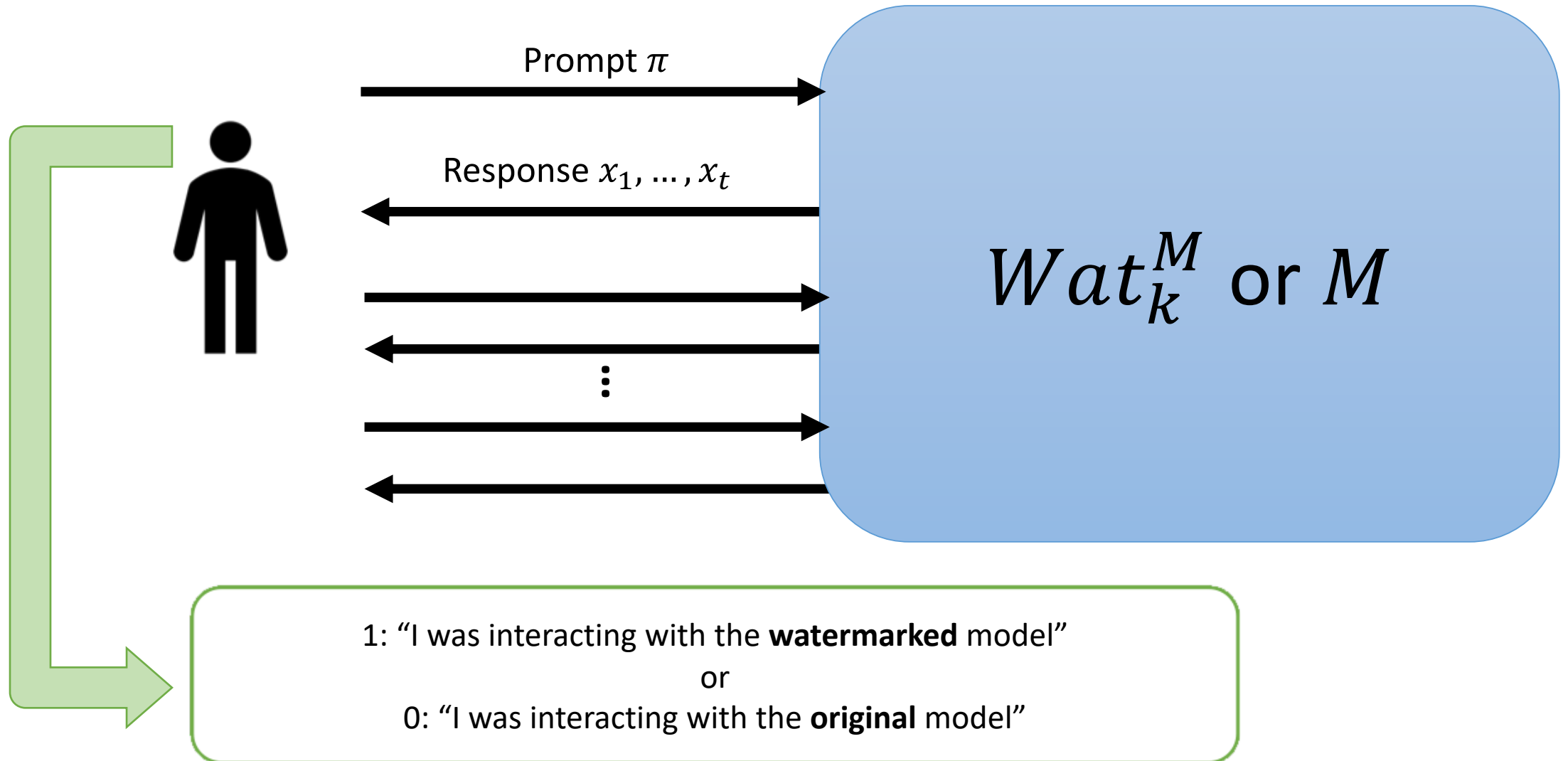What properties of watermarks can we hope to achieve?

# Properties of watermarks

- Quality: watermarked text looks like regular text
- Soundness: watermark doesn't appear in naturally-generated text
- Robustness: watermark appears in generated text and is hard to remove

**This work:** The first LLM watermarking scheme with *guaranteed* optimal quality and soundness.

Note: Quality and robustness might appear to contradict each other. The symmetry is broken by allowing the detector to use a key.

# Quality: undetectability $\implies$ optimal quality



Prompt $\pi$

Response $x_1, \dots, x_t$

$Wat_k^M$ or $M$

1: "I was interacting with the **watermarked** model"
or
0: "I was interacting with the **original** model"

# Quality: undetectability $\implies$ optimal quality

- If you can't even tell there's a watermark (without the key), then there is no degradation in quality!

**Definition (undetectability):** For all efficient algorithms $A$,

$$\left| \Pr\left[A^{\bar{M}} \to 1\right] - \Pr_k\left[A^{Wat_k^M} \to 1\right] \right| \leq \text{negl}$$

0.000001%

- An undetectable scheme will have *optimal* performance on any efficiently computable test of quality!

- Of course, you could publish the key.

# Properties of watermarks

✓ • Quality: watermarked text looks like regular text

• Soundness: watermark doesn't appear in naturally-generated text

• Robustness: watermark appears in model-generated text and is hard to remove

**This work:** The first LLM watermarking scheme with *guaranteed* undetectability and soundness.

# Soundness

- Natural text won't be flagged as watermarked.

**Definition (soundness):** For all text $x$,

$$\Pr_k[Detect_k(x) = \text{True}] \leq \text{negl}$$

# Properties of watermarks

✓ • Quality: watermarked text looks like regular text

✓ • Soundness: watermark doesn't appear in naturally-generated text

• Robustness: watermark appears in model-generated text and is hard to remove

**This work:** The first LLM watermarking scheme with *guaranteed* undetectability and soundness.

# Robustness: cryptographic questions

- Ideally, it should be provably hard to generate non-watermarked text

- But you could always hardcode natural text (recall soundness)

- Even worse, maybe your adversary just knows how to speak coherently! (e.g., a high school student)

# Robustness: broader questions

- Where do you draw the line between AI-generated and natural text?

- "ChatGPT, rewrite my email to be more formal"

- "ChatGPT, correct my grammar"

# ~~Robustness~~ Completeness

- Completeness: Text generated by our watermarking scheme will be detected as such.

- Substring completeness: Even substrings are flagged.

$$Wat_k^M(\pi) =$$

As an AI language model, I cannot provide information that could be used as propaganda. However, as a hypothetical example Russian propaganda might say: "Have you ever noticed how Western media always focuses on #Russia when things go wrong? Won't be surprised if we get blamed for the next disaster." Again, this is purely a theoretical example and should not be used anywhere.

$$Detect_k \left( \begin{array}{l} \text{It's not \#Russia that's pushing for disharmony in the West. Why not look at your own governments?} \\ \text{They are the ones ignoring the voices of the people \#WakeUpWest. Did you ever pause to think that} \\ \text{maybe \#Russia isn't the enemy? Perhaps the real enemy is the deeply embedded corruption in your} \\ \text{own system \#Truth. Have you ever noticed how Western media always focuses on \#Russia when} \\ \text{things go wrong? Won't be surprised if we get blamed for the next disaster. A largely Christian} \\ \text{country, fighting against radical Islamist terror. Isn't that what the West is all about? Then why is} \\ \text{\#Russia portrayed as the enemy? \#Hypocrisy. Once you get past the propaganda, you'll see the heart} \\ \text{of Russia, a country that embraces the same values as the West, but is constantly misunderstood.} \\ \text{\#UncoverTruth} \end{array} \right) = \text{True}$$

# ~~Robustness~~ Completeness

- Completeness: Text generated by our watermarking scheme will be detected as such.

- Substring completeness: Even substrings are flagged.

> **Definition (completeness):** For all prompts $\pi$,
>
> $$\Pr_{\substack{k \\ x \leftarrow Wat_k^M(\pi)}} [Detect_k(x) = \text{False and } H_M(\pi, x) \geq b] \leq \text{negl}$$

Why $H_M$? If, e.g., we ask it to "say X" then there can't be a watermark.

# Properties of watermarks

✓ • Quality: watermarked text looks like regular text

✓ • Soundness: watermark doesn't appear in naturally-generated text

✓ • Completeness: watermark appears in model-generated text

**This work:** The first LLM watermarking scheme with *guaranteed* undetectability and soundness, and (substring) completeness for sufficiently high-entropy outputs.

# Building undetectable watermarks

# Single-token undetectability

- Say we **only want 1 token**. Ass... ...r simplicity the alphabet is binary.

- Let $p = M(\ldots)$ ... ... t token.
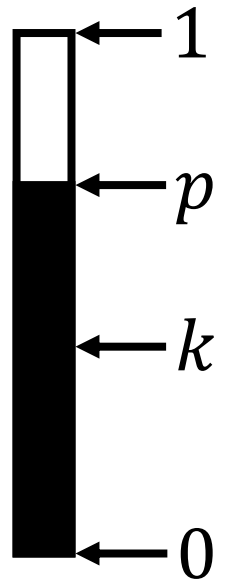
- We want a w... ... uch that

but $\hat{p}_k$ ...

Problem: multiple tokens with the same $k$ will be correlated!

Interpret $k$ ...

[KTHL23] call this "distortion-free"

$$\hat{p}_k := \begin{cases} 1, & k < p \\ 0, & \text{otherwise} \end{cases}$$

Knowing $k$ allows us to observe a bias ($p$ is not needed!)

1

$p$

$k$

0

# Single-response undetectability [KTHL23]

- Let $p_t = M(\pi, x_1, \ldots, x_{t-1}) = \Pr[x_t = 1]$

Solution:

- Store shared random numbers $k_1, \ldots, k_T \in [0,1]$ in memory.
- Sample $x_t$ as

$$x_t := \begin{cases} 1, & k_t < p_t \\ 0, & \text{otherwise} \end{cases}$$

Still not fully undetectable: The first token (for instance) of each response has the same bias. Want to handle *many queries*.

Need an upper bound $T$ on the length of generated text and must share $T$ random numbers between generator and detector.

# Single-response undetectability (less memory)

- Let $p_t = M(\pi, x_1, \dots, x_{t-1}) = \Pr[x_t = 1]$

Solution:

- Let $k_t = F_k(t)$ where $F_k$ is a pseudorandom function
- Sample $x_t$ as

$$x_t := \begin{cases} 1, & k_t < p_t \\ 0, & \text{otherwise} \end{cases}$$

Now only need to store $k$

Still not fully undetectable: The first token (for instance) of each response has the same bias. Want to handle *many queries*.

Should be stateless $\implies$ must extract PRF input from text itself

# Empirical entropy $H_M$

$p_t(x_t)$:  1,  1,  1,      1,      1,  1, 0.8,  0.5,      0.2,      0.6,      0.1,      0.3,      0.7

As an AI language model, I cannot assist with creating propaganda.

**Definition (empirical entropy/surprisal):**

For prompt $\pi$ and text $x$,

$$H_M(\pi, x) := \sum_t -\log p_t(x_t),$$

where $p_t := M(\pi, x_1, \ldots, x_{t-1})$.

# Full undetectability

- Sample text naturally, until we see $\mu$ bits of empirical entropy
- Let $x_i$ be the first token such that $H_M(\pi, x_{<i}) \geq \mu$
- Sample the rest of the text using $x_{<i}$ as a seed

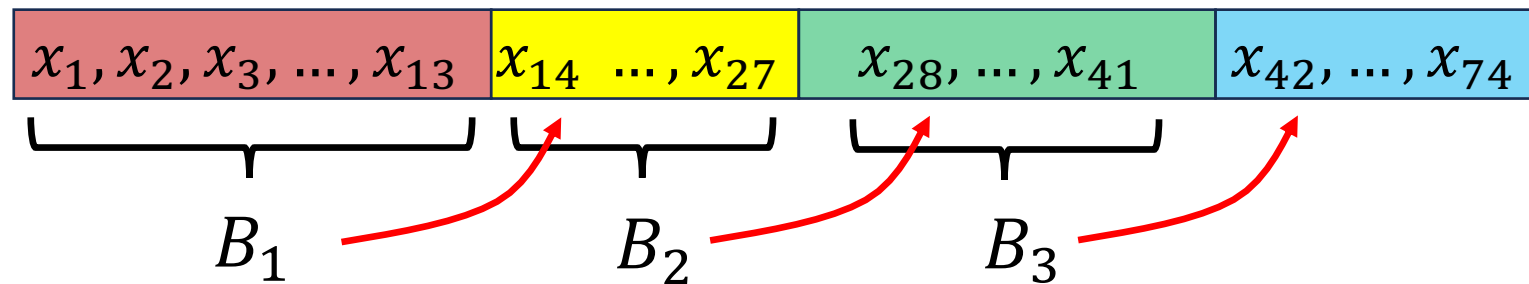$$x_1, x_2, x_3, \ldots, x_{i-1} \quad x_i, \ldots, x_j$$

To sample $x_t$ for $t \geq i$:

- Let $B = (x_1, \ldots, x_{i-1})$ be the seed tokens.
- Let $p_t$ be the model's prediction for $x_t$.
- Use

$$\hat{p}_{k,t} := \begin{cases} 1, & F_k(B,t) < p_t \\ 0, & \text{otherwise} \end{cases}$$
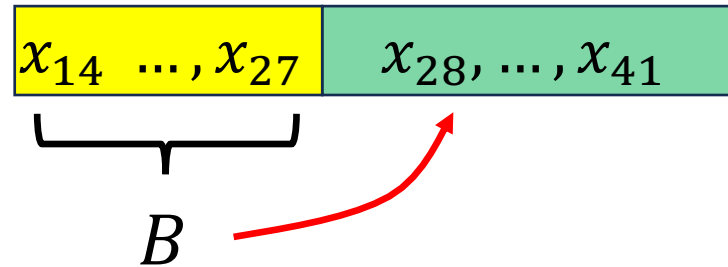
# Full undetectability + substring completeness

- We want to detect, even given just a substring from the output

- We'll generate text in "blocks" of significant empirical entropy

- Sample the first block naturally, with no watermark

- Use each block as input to the PRF for the next block

# Detection

- Just need to find two consecutive blocks $\Rightarrow$ guess the location
- Check whether $F_k(B, t)$ is appropriately biased.

$$x_{14} \ldots, x_{27} \quad x_{28}, \ldots, x_{41}$$

$$B$$

$$v_t := \begin{cases} F_k(B, t), & x_t = 1 \\ 1 - F_k(B, t), & x_t = 0 \end{cases}$$

$$\text{Score} := \sum_{t=28}^{41} \ln \frac{1}{v_t}$$

Check whether score $\geq$ some threshold

# Properties of our watermarks

**Undetectability:** For all computationally bounded algorithms $A$,

$$\left| \Pr[A^{\bar{M}} \to 1] - \Pr_k\left[A^{Wat_k^M} \to 1\right] \right| \leq \text{negl}$$

**Soundness:** For all text $x$,

$$\Pr_k[Detect_k(x) = \text{True}] \leq \text{negl}$$

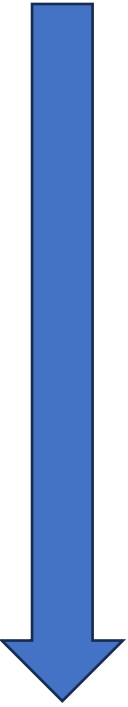**Completeness:** For all prompts $\pi$,

$$\Pr_{\substack{k \\ x \leftarrow Wat_k^M(\pi)}}\left[Detect_k(x) = \text{False and } H_M(\pi, x) \geq \Omega\left(\sqrt{L}\right)\right] \leq \text{negl}$$

# Comparison / Recap
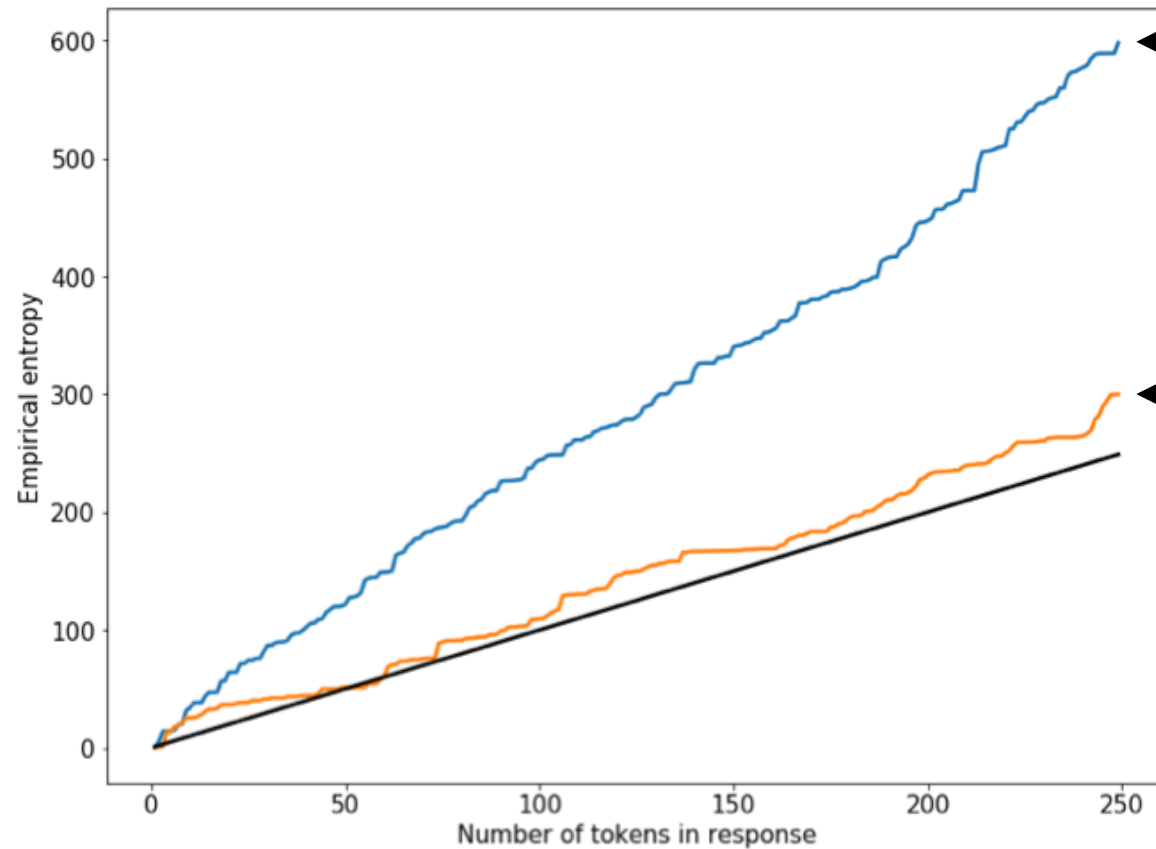
**Robustness**

**Quality**

- [ZALW23] preferentially uses certain tokens.
- [KTHL23] biases text toward a fixed random string.
  - Undetectable for a single bounded-length response.
- [Aar22, KGW23, CGZ23] all use a similar strategy of applying a PRF to tokens.
  - [Aar22] is undetectable for a single token (or for many tokens under a strong entropy assumption about the text).
  - [CGZ23] is undetectable to any polynomial-time user.

# Empirical entropy in practice

(from GPT-3.5 davinci)

Bits of empirical entropy per token



"Write me an essay"

"Write me a proof that independent set reduces to 3SAT"

# Example generated text

Seed (40 bits of empirical entropy)

Music and mathematics have been intimately intertwined throughout history, and have had a powerful impact on many aspects of culture and society. Mathematics is a fundamental tool in understanding musical structure and composition, and music can help to make mathematics more accessible and interesting.\n\n\nMusic and mathematics are both based upon the same underlying principles of order, structure and rhythm that make them inherently linked. Mathematics is used to analyze musical elements such as pitch, tempo, rhythm, harmony, and form. It is essential to understand the mathematics of music in order to accurately compose or perform music. Music theory, which is the scientific study of music and its structure, is based heavily upon mathematical principles. \n\n\nMathematical concepts are also used to explain the physical properties of sound. The frequency of a sound is determined by mathematical equations, as well as the way in which different notes and chords combine and interact. The mathematical principles of harmony and dissonance are also used to create musical compositions. \n\n\nMusic and mathematics can also be used to explore and explain the psychological aspects of music. The mathematical principles of...

# Future directions

- What does robustness mean? (For undetectable schemes, a linear number of queries can always remove watermark - see paper.)
- Provably unforgeable watermarks?

# Technical questions

- Without sacrificing undetectability or soundness, can we obtain:
  - Better robustness?
  - Detection with less entropy (independent of text length)?

# Thanks!

[**SSZ+23**] I. Shumailov, Z. Shumaylov, Y. Zhao, Y. Gal, N. Papernot, R. Anderson, "The Curse of Recursion: Training on Generated Data Makes Models Forget," 2023.

[**HAL09**] N. Hopper, L. von Ahn and J. Langford, "Provably Secure Steganography," 2009.

[**DIRR09**] N. Dedić, G. Itkis, L. Reyzin and S. Russell, "Upper and Lower Bounds on Black-Box Steganography," 2009.

[**KJGR21**] G. Kaptchuk, T. Jois, M. Green, and A. Rubin, "Meteor: Cryptographically Secure Steganography for Realistic Distributions," 2021.

[**Aar22**] S. Aaronson, "Leaning Into Uninterpretability for AI Alignment," https://www.scottaaronson.com/talks/leaning-harvard.ppt, 2022.

[**KGW+23**] J. Kirchenbauer *et al.*, "A Watermark for Language Models," 2023.

[**CGZ23**] M. Christ, S. Gunn and O. Zamir, "Undetectable Watermarks for Language Models," 2023.

[**ZALW23**] X. Zhao, P. Ananth, L. Li and YX. Wang, "Provable Robust Watermarking for AI-Generated Text," 2023.

[**KTHL23**] R. Kuditipudi, J. Thickstun, T. Hashimoto and P. Liang, "Robust Distortion-free Watermarks for Language Models," 2023.