# Watermarking Makes Language Models Radioactive

December 2024 - NeurIPS (Spotlight) ✨

Tom Sander[1,2], Pierre Fernandez[1,3]                    (equal contributors)

Matthijs Douze[1], Alain Durmus[2], Teddy Furon[3]        (advisors)

[1] Meta, FAIR        [2] Ecole polytechnique        [3] Inria Rennes

∞ Meta

# Motivation

LLM **post-training**

- ○ Requires a lot of high quality annotations and tricks → 🧠 and 💰
- ○ Practitioners train on data output by a model (e.g., GPT4) → **IP issues**



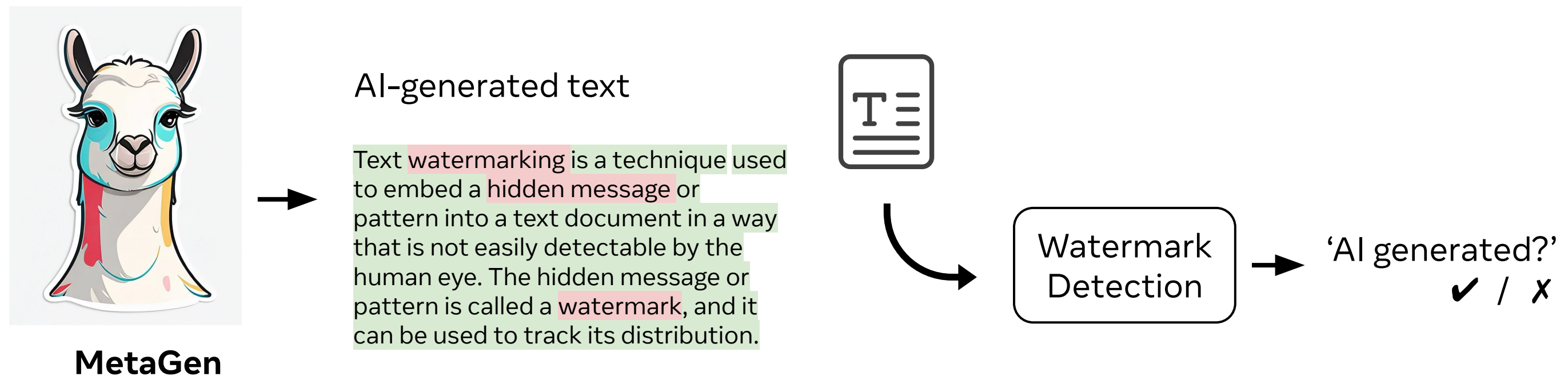Estimated Training Costs of Large Models

# Detection Problem

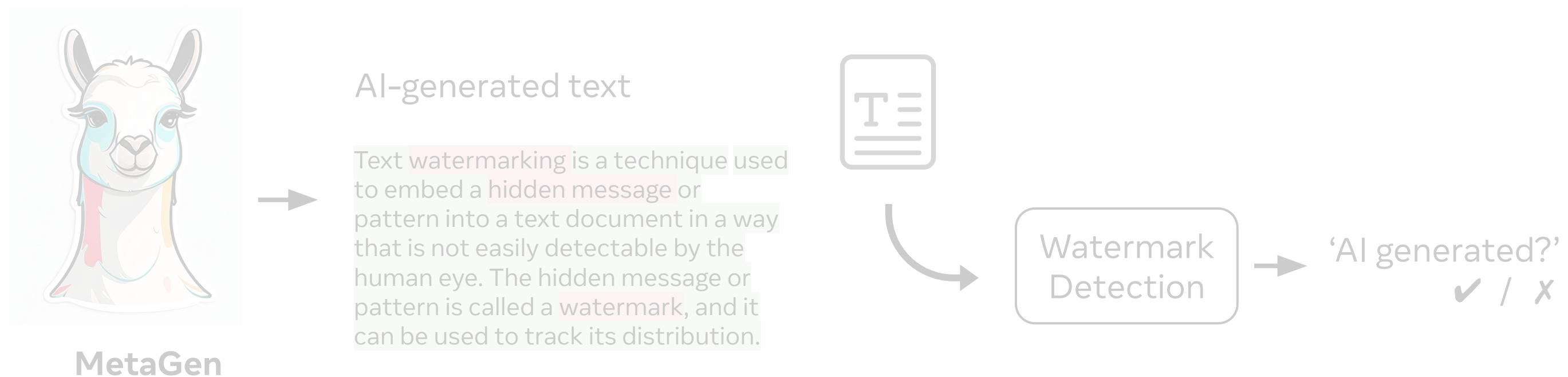"Did Bob train on outputs from Alice's model?" is a **very difficult question**

# Could Watermarking Give the Answer?

- **Watermarking LLMs** outputs ≈ free lunch
  - Keeps **quality** of the generated text
  - Greatly improves **detection**



**MetaGen**

AI-generated text

Text watermarking is a technique used to embed a hidden message or pattern into a text document in a way that is not easily detectable by the human eye. The hidden message or pattern is called a watermark, and it can be used to track its distribution.

Watermark Detection → 'AI generated?' ✔ / ✗

# Could Watermarking Give the Answer?

- **Watermarking LLMs** outputs ≈ free lunch
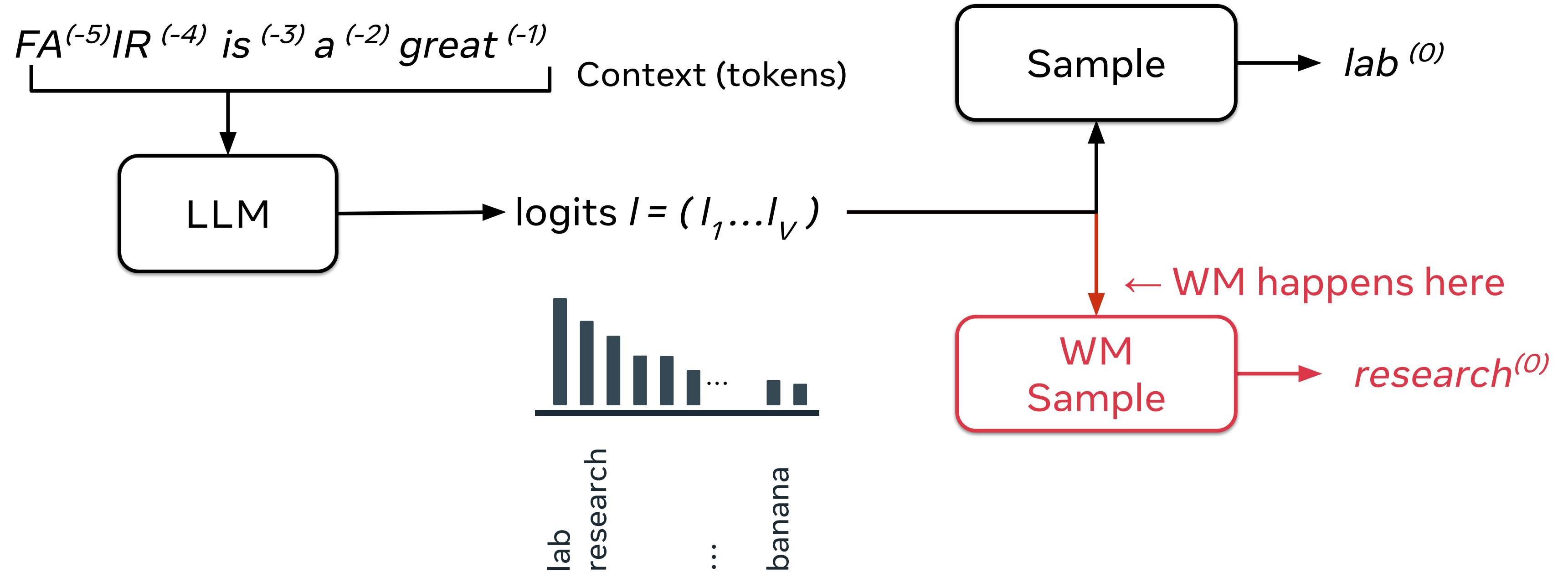  - Keeps **quality** of the generated text
  - Greatly improves **detection**



**MetaGen**

AI-generated text

Text watermarking is a technique used to embed a hidden message or pattern into a text document in a way that is not easily detectable by the human eye. The hidden message or pattern is called a watermark, and it can be used to track its distribution.

Watermark Detection

'AI generated?'
✔ / ✗

→ *"What occurs when we fine-tune an LLM on watermarked data?"*

# LLM Watermarking 101

# Watermarking for LLMs

**Generation with LLMs**

$FA^{(-5)} IR^{(-4)}$ is $^{(-3)}$ a $^{(-2)}$ great $^{(-1)}$    Context (tokens)

LLM → logits $l = (l_1 ... l_V)$

Sample → lab $^{(0)}$

← WM happens here

WM Sample → research $^{(0)}$

lab   research   ...   banana

# First Example – Kirchenbauer et al.

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

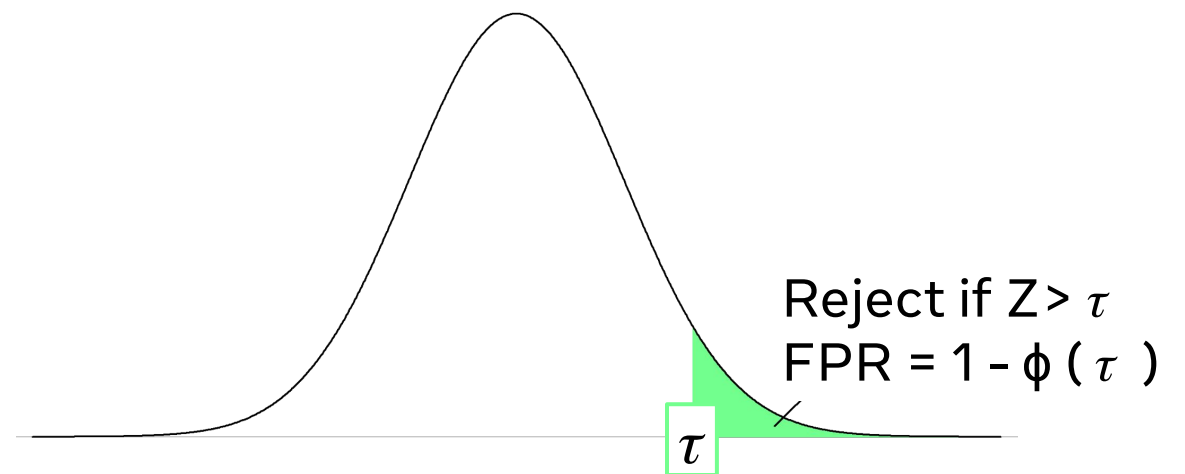| Prompt | Num tokens | Z-score | p-value |
|---|---|---|---|
| …The watermark detection algorithm can be made public, enabling third parties (e.g., social media platforms) to run it themselves, or it can be kept private and run behind an API.  We seek a watermark with the following properties: | | | |
| **No watermark**<br>Extremely efficient on average term lengths and word frequencies on synthetic, microamount text (as little as 25 words)<br>Very small and low-resource key/hash (e.g., 140 bits per key is sufficient for 99.999999999% of the Synthetic Internet | 56 | .31 | .38 |
| **With watermark**<br>- minimal marginal probability for a detection attempt.<br>- Good speech frequency and energy rate reduction.<br>- messages indiscernible to humans.<br>- easy for humans to verify. | 36 | 7.4 | 6e-14 |

# Count Greenlist/Redlist Tokens



Without Watermark

Research on human aggression has been going on for decades , and has been done in a variety of ways . There is no single way to answer this question , and it will depend on the research method and the research question . Humans go to war to protect or expand their territories , and this often results in violence and death . Some argue that war is a natural part of humanity and cannot be eliminated , but this is a controversial view . Human violence has not always

Green: 44 / 91

**Human/Unwater marked**

With Watermark

To see if they have the stomach to kill another human , and to see if they have the stomach to kill another living thing and risk their life . I 'm not saying that justifies it , but those are the reasons . " They had the stomach to kill another human " That 's not very nice . They have to kill to survive , so they have to kill someone to kill another person They don 't HAVE TO kill to survive . There are other methods of getting food other than killing . Humans invented agriculture

Green: 73 / 99

**Watermarked**

## Statistical test

- Total score $S = \sum_{t \in 1,..,T} S_t$ = number of greenlist tokens
- $H_0$ = "text is not watermarked"
- Reject based $H_0$ on Z-Score:  $Z = (S-\mu)/\sigma$  $(= (S-T/2)/T/4)$

Reject if $Z > \tau$
FPR = $1 - \phi(\tau)$

$\tau$

9

# How to Choose Greenlist/Redlist?

✗ Fixed lists

→ heavily biases the generation

⇔ "Generate a text on France without using the word France"

lab 🟧

research 🟩

France 🟧

water 🟧

mark 🟩

FA 🟩

IR 🟧

word 🟧

…

# How to Choose Greenlist/Redlist?

✗ Fixed lists

→ heavily biases the generation

✔ Make it dependant on previous tokens

"After word 'water', greenlist/redlist are …"

# Sampling with Greenlist/Redlist

$FA^{(-5)}IR^{(-4)}$ $is^{(-3)}$ $a^{(-2)}$ $great^{(-1)}$

Context (tokens)

LLM

lab
research
place
French
...

# Sampling with Greenlist/Redlist

Get greenlist / redlist

↑

previous token(s) /
watermark window

$FA^{(-5)}IR^{(-4)}$ $is^{(-3)}$ $a^{(-2)}$ $great^{(-1)}$

Context (tokens)

LLM

lab
research
place
French
...

# Sampling with Greenlist/Redlist

Get greenlist / redlist

previous token(s)/
watermark window

+ δ=1.0 to greenlist tokens' logit

$FA^{(-5)}IR^{(-4)}$ is $^{(-3)}$ a $^{(-2)}$ great $^{(-1)}$

Context (tokens)

LLM

lab
research
place
French
...

*softmax*

$\longrightarrow p$

# Sampling with Greenlist/Redlist

Get greenlist / redlist

previous token(s) /
watermark window

+ δ=1.0 to greenlist tokens' logit

$FA^{(-5)} IR^{(-4)} is^{(-3)} a^{(-2)} great^{(-1)}$

Context (tokens)

LLM

lab
research
place
French
...

softmax

$p$

Multinomial
sampling

research

# Detection with Greenlist/Redlist

Compute score

Get greenlist / redlist

previous token(s) /
watermark window

$FA^{(-5)}IR^{(-4)}\ is^{(-3)}\ a^{(-2)}\ great^{(-1)}\ research^{(0)}$

S += 1

Statistical test

- Total score $S = \sum_{t \in 1,..,T} S_t$ = number of greenlist tokens
- $H_0$ = "text is not watermarked"
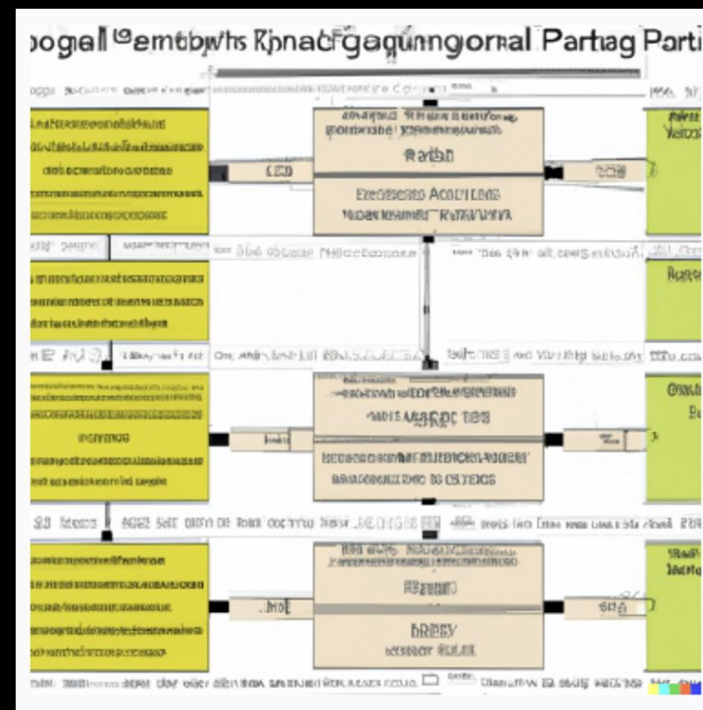- Reject based $H_0$ on Z-Score:  $Z = (S-\mu)/\sigma$  or Binomial test

# Second Example - Aaronson et al.

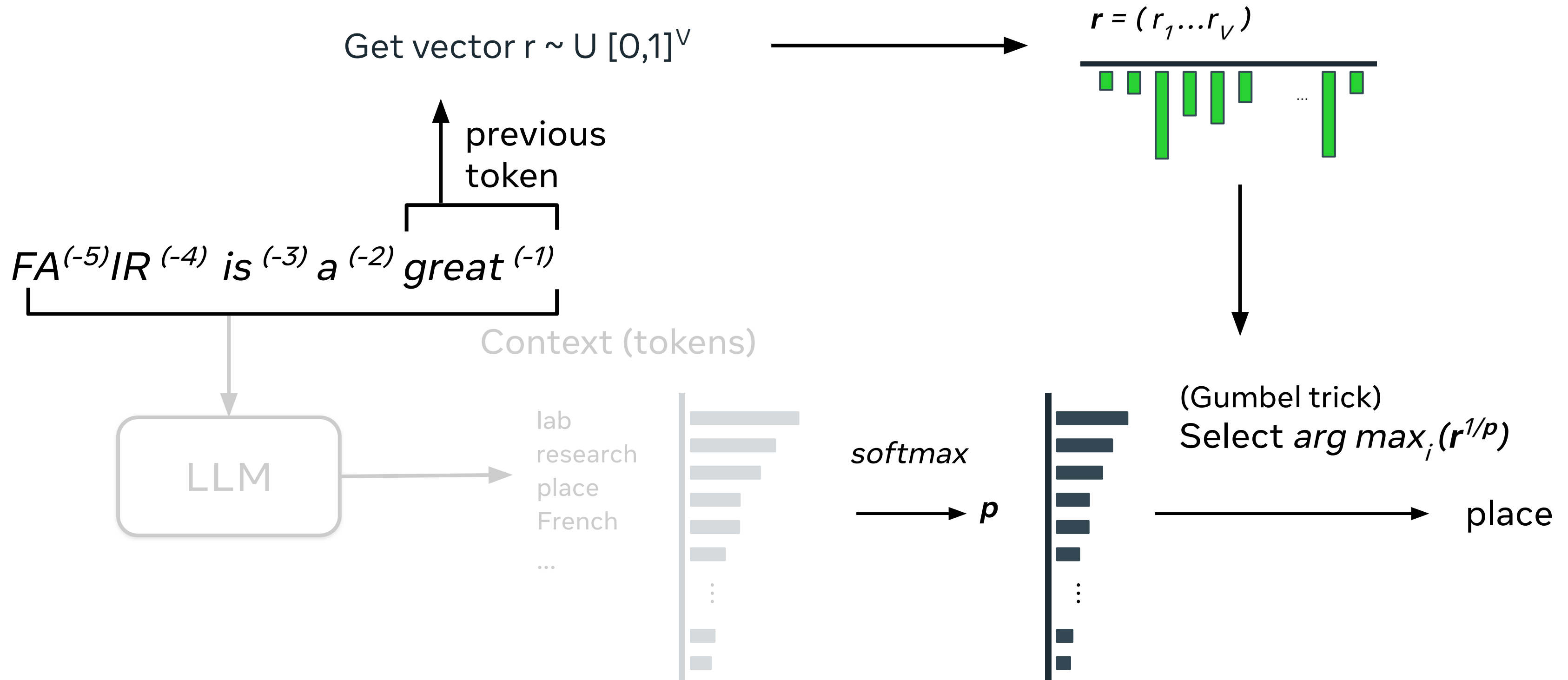📄 Aaronson et al., *Watermarking GPT Outputs*, 2022

# Sampling with Gumbel Trick

Get vector r ~ U $[0,1]^V$

$r = ( r_1 ... r_V )$



previous token

$FA^{(-5)}IR^{(-4)}$ is $^{(-3)}$ a $^{(-2)}$ great $^{(-1)}$

Context (tokens)

LLM

lab
research
place
French
...

*softmax*

$p$
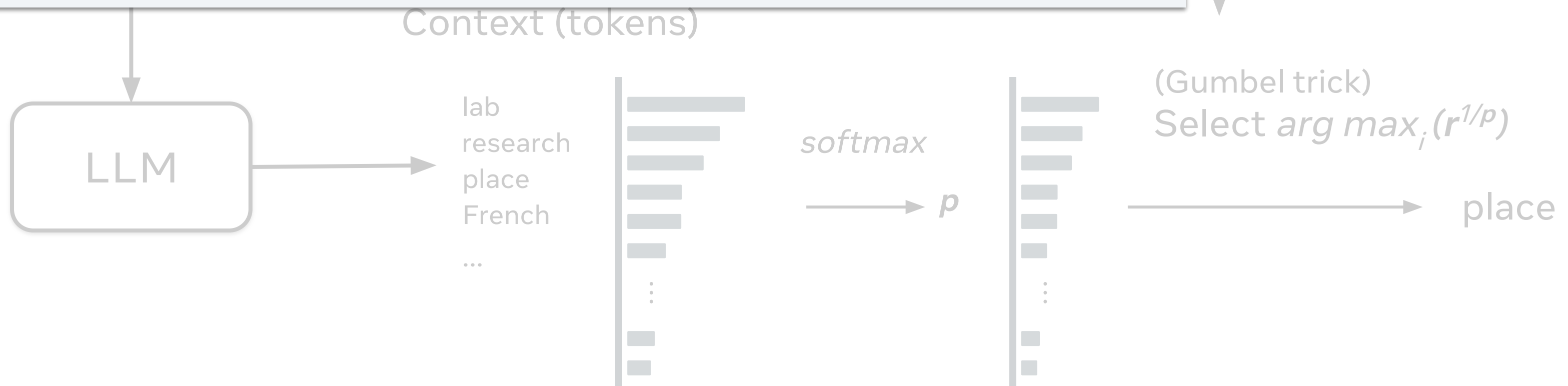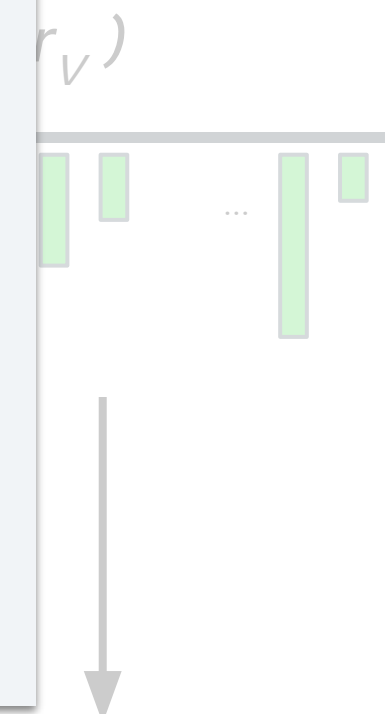
(Gumbel trick)
Select *arg max$_i$ ($r^{1/p}$)*

place

# Sampling with Gumbel Trick

**Property (Gumbel trick):**

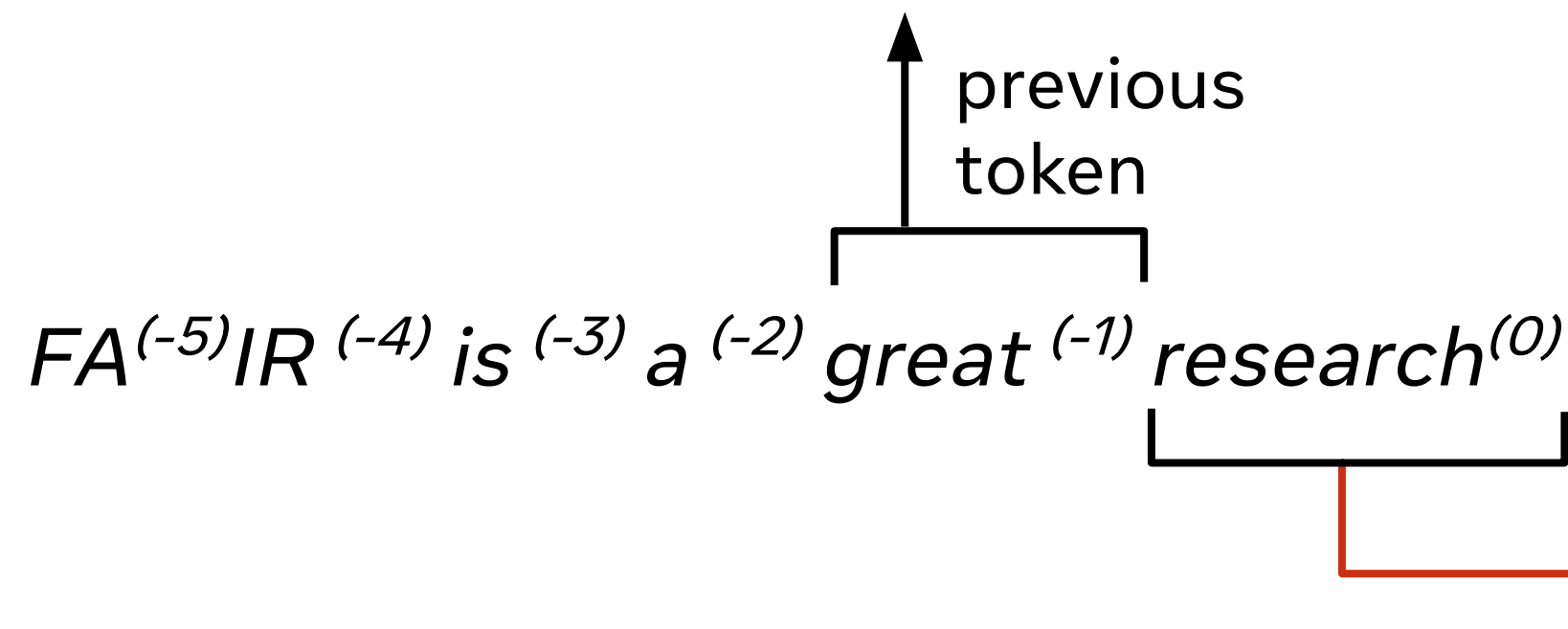$$\forall i \in [1, .., V],\ P\left(\arg\max_i R_i^{1/p_i} = i \mid R \sim U(0,1)^V\right) = p_i$$

$\rightarrow$ **"Proba of choosing token $i$ is $p_i$"**

Context (tokens)

lab
research
place
French
...

*softmax*

$p$

(Gumbel trick)
Select $arg\ max_i(r^{1/p})$

place

LLM

# Detection with Z-score

Compute score

Get vector $r \sim U[0,1]^V$

previous
token

$FA^{(-5)}IR^{(-4)} is^{(-3)} a^{(-2)} great^{(-1)} research^{(0)}$

Score increment:
$-\ln(1-r_t)$ with t the index of chosen token

$S +\!= -\ln(1-r_t)$

Statistical test

- Total score $S = \sum_{t \in 1,..,T} S_t$
- $H_0$ = "text is genuine"   –   $H_1$ = "text is watermarked"
- Reject based $H_0$ on Z-Score: $Z = (S-\mu)/\sigma$

Reject if $Z > \tau$
$FPR = 1 - \phi(\tau)$

$\tau$

# Detection with Z-score

Compute score

Get vector r ~ U [0,1]$^V$

previous

$FA^{(\ldots}$

osen token

**Property:**

( H$_0$ ) For non-watermarked texts: $\quad \mathbb{E}(S_T) = T$

( H$_1$ ) For watermarked texts: $\quad \mathbb{E}(S_T) \geq T + \left( \dfrac{\pi^2}{6} - 1 \right) H_T$

Sta

- Total score S = $\sum_{t \in 1,\ldots,T} S_t$
- H$_0$ = "text is genuine"  –  H$_1$ = "text is watermarked"
- Reject based H$_0$ on Z-Score:  Z = (S-μ)/σ

Reject if Z > $\tau$

FPR = 1 - ϕ ( $\tau$ )

$\tau$

# Detection with Z-score

Compute score

Get vector r ~ U $[0,1]^V$

previous
token

$FA^{(-5)}IR^{(-4)}\ is^{(-3)}\ a^{(-2)}\ great^{(-1)}\ research^{(0)}$

Score increment:
$-\ln(1-r_t)$ with t the index of chosen token

$S \mathrel{+}= -\ln(1-r_t)$

## Statistical test

- Total score $S = \sum_{t \in 1,..,T} S_t$
- $H_0$ = "text is genuine"  –  $H_1$ = "text is watermarked"
- Reject based $H_0$ on Z-Score: $Z = (S-\mu)/\sigma$
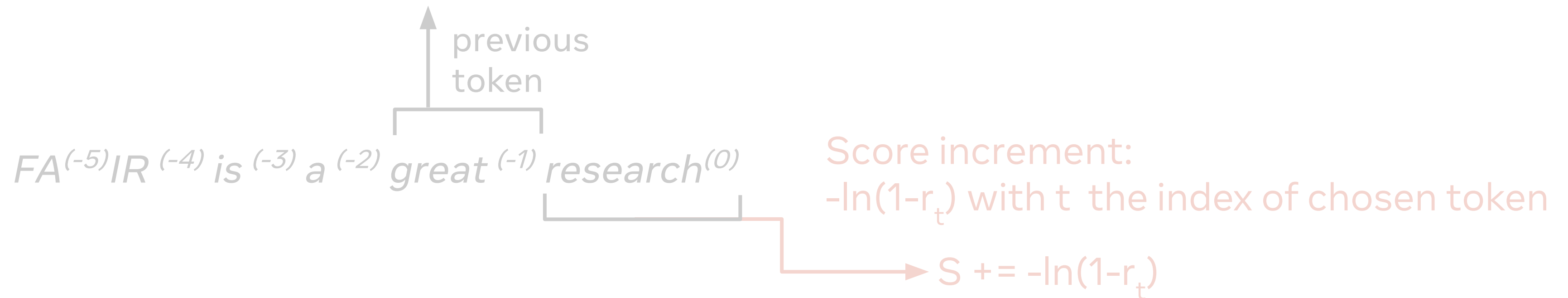
Reject if Z > $\tau$
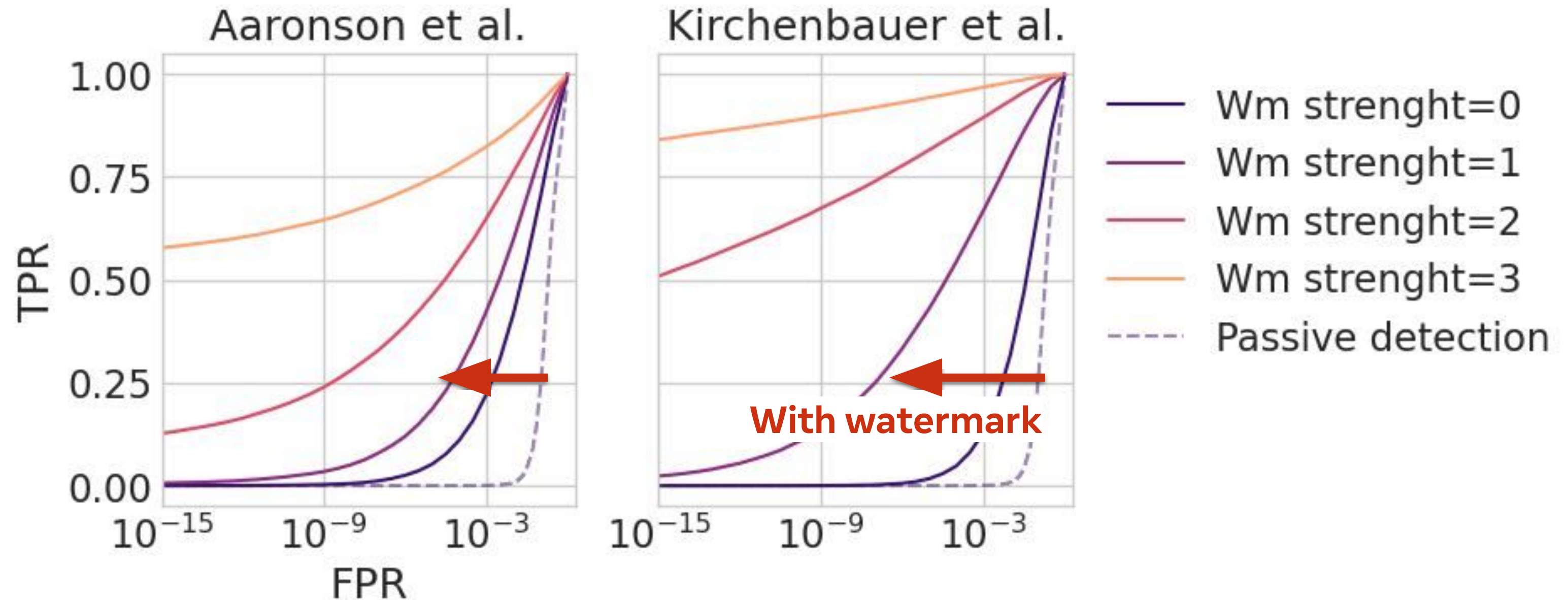FPR = 1 - φ ($\tau$)

$\tau$

# Example - Detection Results

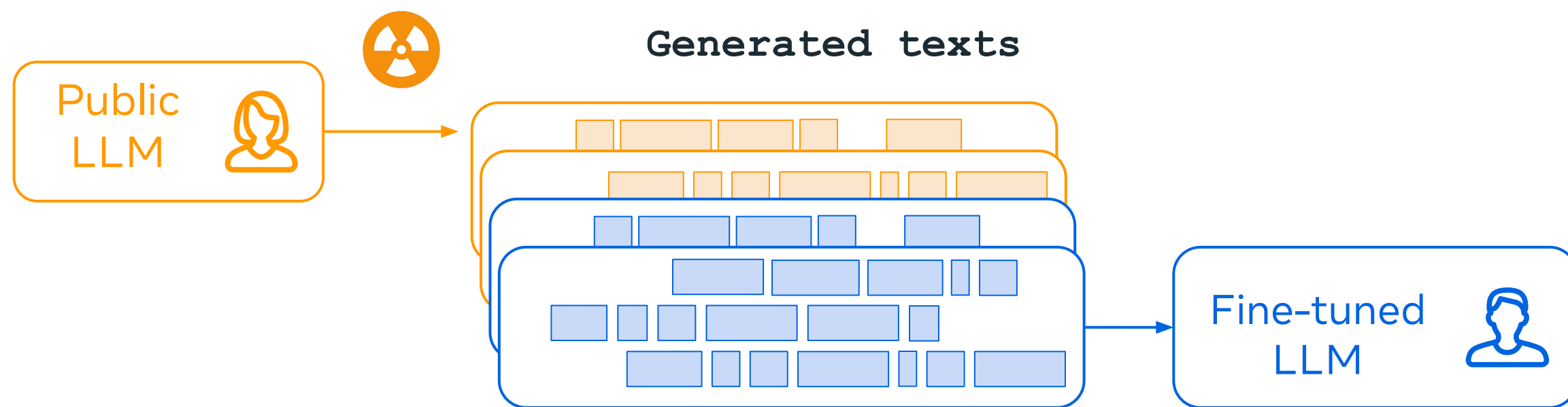10k positive AI-generated texts (from OpenAssistant Conversations dataset)

Passive detection ↔ DetectGPT [📄 Mitchell, Eric, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. "Detectgpt: Zero-shot machine-generated text detection using probability curvature.", ICML 2023]
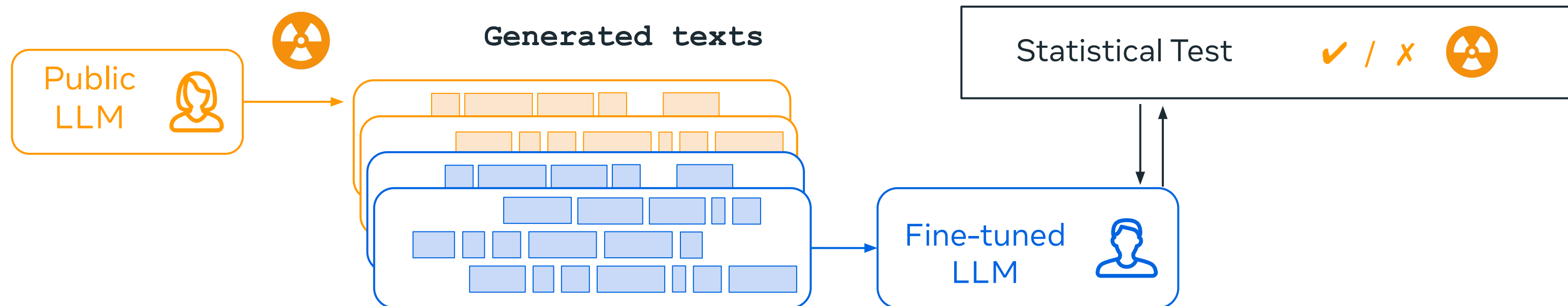
# Radioactivity

# Problem under Study

- Bob **fine-tunes** his LLM on **training data** with a small proportion of texts **coming from Alice's LLM.**

# Problem under Study

- Bob **fine-tunes** his LLM on **training data** with a small proportion of texts **coming from Alice's LLM.**

- Alice wants to know if Bob has **fine-tuned on outputs from her model**

# Radioactivity

**Definition:** Radioactivity refers to the possibility for Alice to detect with statistical evidence that Bob fine-tuned on outputs from her model

More rigorously,

**Definition 1** (Text Radioactivity). *Dataset $D$ is $\alpha$-radioactive for a statistical test $T$ if "$\mathcal{B}$ was not trained on $D$" $\subset \mathcal{H}_0$ and $T$ is able to reject $\mathcal{H}_0$ at a significance level ($p$-value) smaller than $\alpha$.*

**Definition 2** (Model Radioactivity). *Model $\mathcal{A}$ is $\alpha$-radioactive for a statistical test $T$ if "$\mathcal{B}$ was not trained on outputs of $\mathcal{A}$" $\subset \mathcal{H}_0$ and $T$ is able to reject $\mathcal{H}_0$ at a significance level smaller than $\alpha$.*
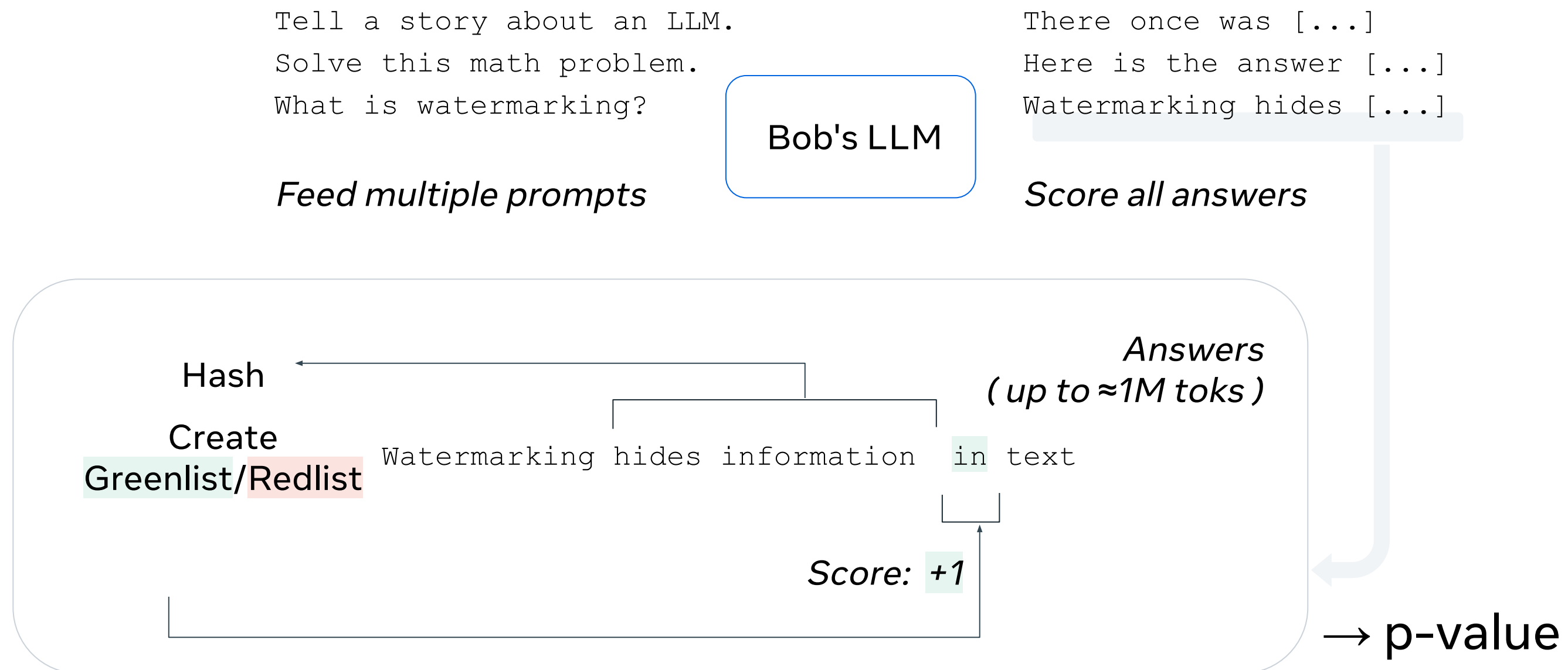
# Different Settings

**Model access**

| | Model is open<br>(Mistral, Llama, Gemma, etc.) | API access only<br>(GPT, Claude, etc.) |
|---|---|---|
| **Access to the text used by Bob**<br>(GPT, Claude, etc.) | Open / Supervised | Closed / Supervised |
| **Text used by Bob is unknown**<br>(Llama, API but obfuscation of user) | Open/ Unsupervised | Closed/ Unsupervised |

**Data access**

Radioactivity detection availability from other methods in the literature

| | With WM | | Without WM (MIA) | | IPP | |
|---|---|---|---|---|---|---|
| | Open | Closed | Open | Closed | Open | Closed |
| Supervised | ✓ | ✓ | ✓ | ✗ | ✓ | ∼ |
| Unsupervised | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |

# Naive Approach for Radioactivity Detection with Watermarking

Prompt the model, get many output tokens, get the score and the p-value of the WM detection

```
Tell a story about an LLM.
Solve this math problem.
What is watermarking?
```

Bob's LLM

```
There once was [...]
Here is the answer [...]
Watermarking hides [...]
```

*Feed multiple prompts*

*Score all answers*

Hash

Create
Greenlist/Redlist

*Answers
( up to ≈1M toks )*

Watermarking hides information  in text

*Score: +1*

→ p-value

# Problems with the Naive Approach

**Watermark signal is <u>weak</u>**

→ hard to get p-values < $10^{-1}$ for low proportions of watermarked data in the training set

**p-values break down when too many tokens are scored**

→ when scoring two many tokens, the detection test gives very low p-values even for LLMs trained without watermarked text, so the statistical tests are inaccurate

# Problems with the Naive Approach

**Watermark signal is <u>weak</u>**

→ hard to get p-values < $10^{-1}$ for low proportions of watermarked data in the training set

**p-values break down when too many tokens are scored**

→ when scoring two many tokens, the detection test gives very low p-values even for LLMs trained without watermarked text, so the statistical tests are inaccurate
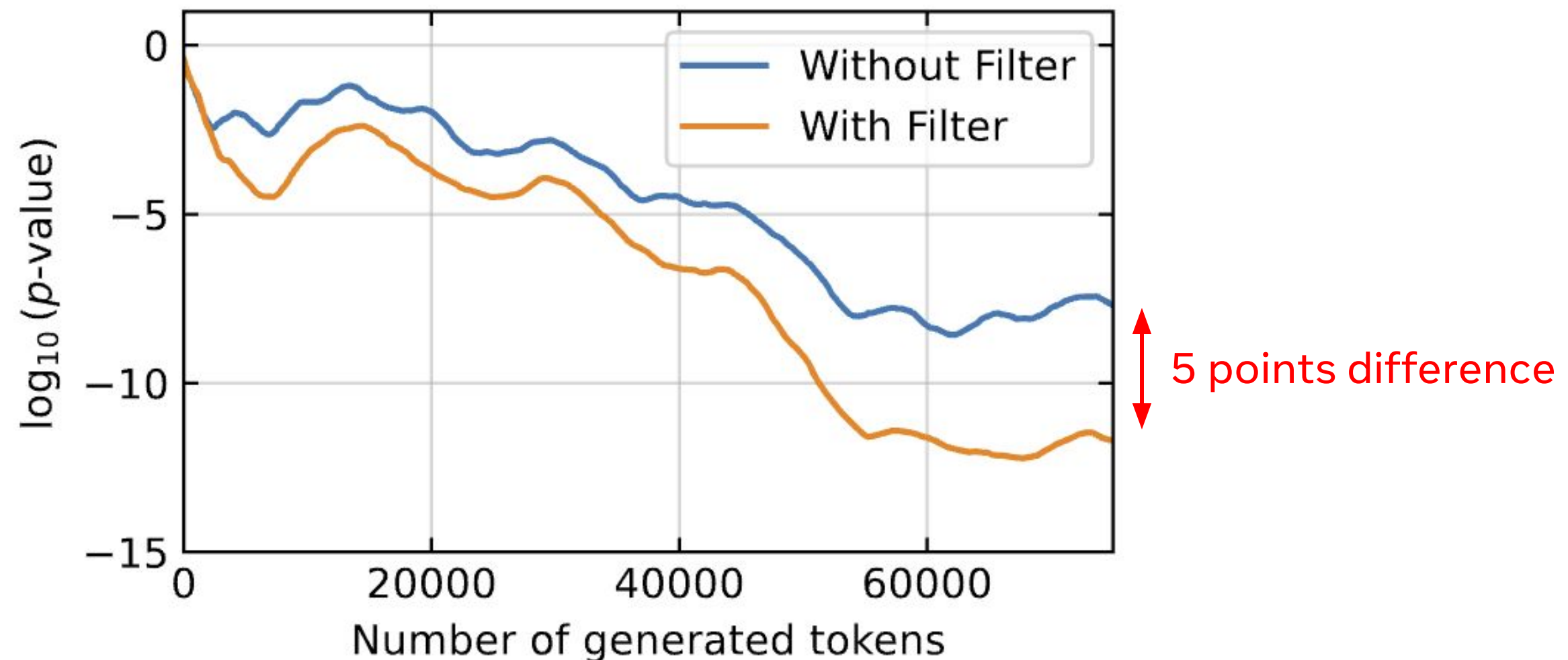
**Improvements**

- Leverage access to the data
- Leverage access to the model
- While keeping accurate p-values through deduplication

# Trick 1: Filter

Radioactivity can only be detected on watermark windows present in training

- Supervised setting: only score watermark windows suspected to be part of training
- Unsupervised setting: see what are the watermark windows that are most often produced by the watermark, and only score these
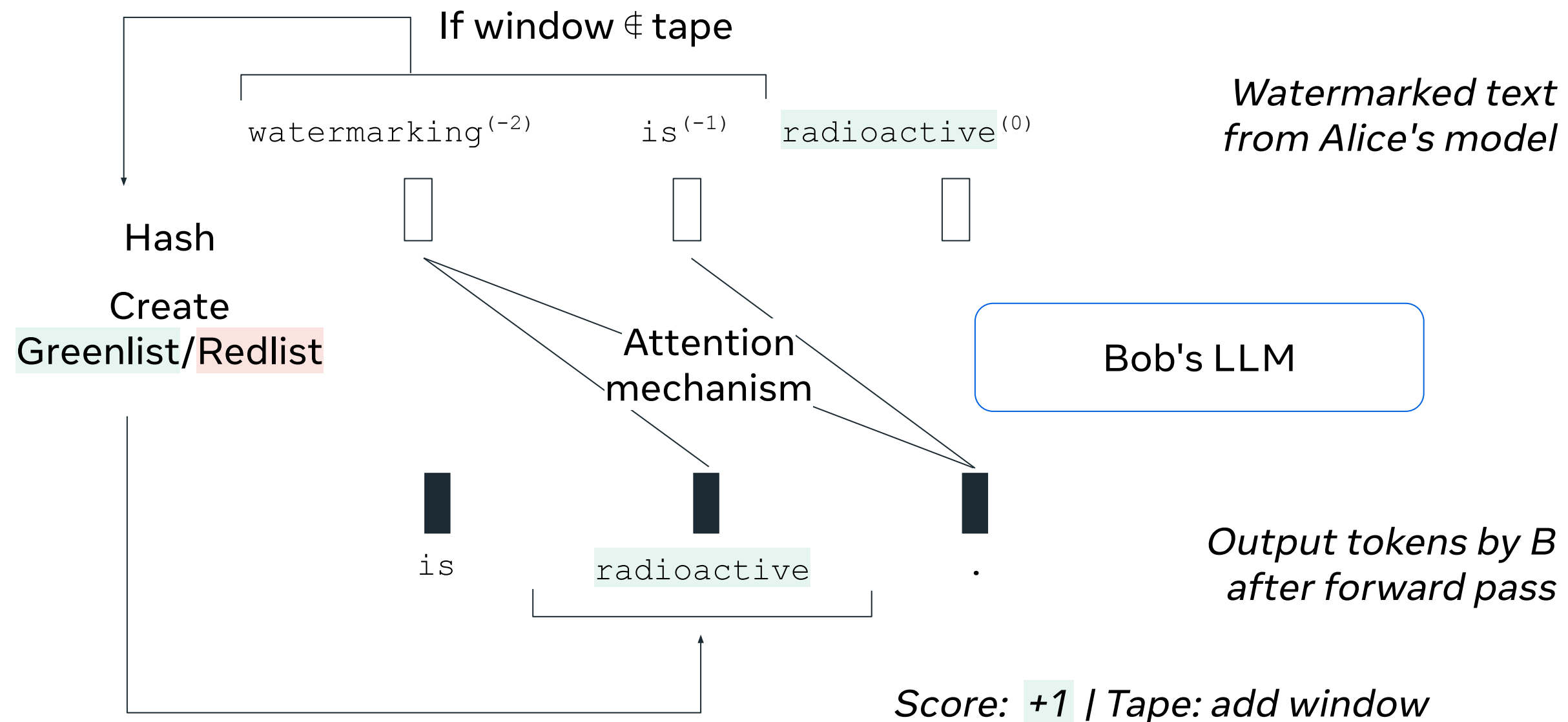
# Trick 2: Choose the Good Input

Radioactivity can only be detected on k-grams that were present in training

- **Closed-model**: Alice prompts Bob's model with questions that she thinks were used
- **Open-model**: Alice "reads" the data that she thinks Bob has used

# Trick 3: Open Model

When access to the model is given, Alice can forward text directly to the model

- **Gain in efficiency**: one pass forward only
- **Gain in supervision**: the model sees exact reproduction of watermark window & context

If window $\notin$ tape

$watermarking^{(-2)}$   $is^{(-1)}$   $radioactive^{(0)}$

*Watermarked text from Alice's model*

Hash

Create
Greenlist/Redlist

Attention
mechanism

Bob's LLM

is   radioactive   .

*Output tokens by B after forward pass*

*Score:  +1  | Tape: add window*

# Trick 4: Deduplication for False Positives

- Very important to get reliable p-values

| Access to Model | De-duplication | |
| --- | --- | --- |
| | With | Without |
| Open | $0.46_{\pm 0.27}$ | $0.053_{\pm 0.12}$ |
| Closed | $0.42_{\pm 0.30}$ | $< 10^{-30}$ |

- Lots of rules:
  - Don't score tokens whose watermark window have already been scored
  - Don't score tokens whose watermark window is already in the attention span

# Experimental Setup

1. Generate watermarked instructions with Llama-2-chat-7b and Self-Instruct
2. Fine-tune Llama-1-7b with varying proportions of watermarked instructions
3. Get p-values of radioactivity detection

# Detection Results under the Different Settings



- if the suspect model is open-weight, detection has p-value $< 10^{-5}$ even when as little as 5% of training text is watermarked

- when Alice only has API access but knows which data have been used, detection has p-value $< 10^{-10}$ even when 1% of the training text is watermarked
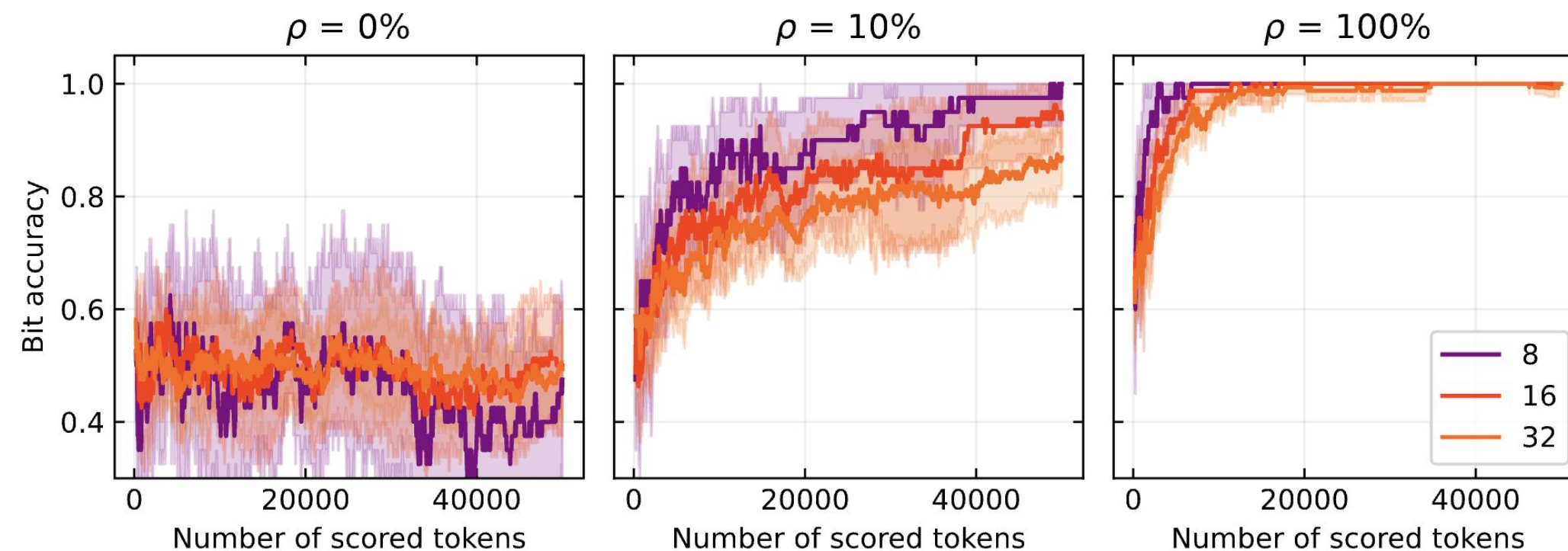
# Ablations

## Post-training optimization has a big influence on radioactivity

$\log_{10}$ p-value for 10k observed tokens under the supervised-open model setting

| (a) Learning rate. | | | | (b) Epoch. | | | | | (c) Adapters. | | | (d) Model size. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $10^{-5}$ | $5 \cdot 10^{-5}$ | $10^{-4}$ | | 1 | 2 | 3 | 4 | | Full | Q-LoRA | | 7B | 13B |
| -32.4 | -49.6 | -58.0 | | -20.8 | -29.2 | -33.2 | -34.8 | | -32.4 | -11.0 | | -32.4 | -33.2 |

## The method generalizes to multi-bit watermarking

# Ablations

**A lot more in the paper!**

# Main Takeaways

**Watermarking** makes LLM radioactive:

- Training on watermarked data can be **detected with very high confidence**…
- … even for **small proportions** of WM data

**Thanks!**