Information-Theoretic Driven Watermarking of Large-Language Models

Dor Tsur, October 17th 2025.







Collaborators

Carol Long Harvard University



Claudio Mayrink Verdun
Harvard University



Hsiang Hsu JP Morgan & Chase



Chun-Fu (Richard) Chen
JP Morgan & Chase



Haim H. Permuter Ben-Gurion University



Sajani Vithana Harvard University



Flavio P. Calmon
Harvard University





- LLMs demonstrate human-level text generation capabilities
- Prone to misuse
 - Misinformation & Fake news
 - Al-generated scams
 - Academic dishonesty

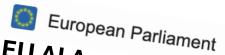


Malwarebytes

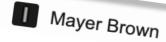
AI-supported spear phishing fools more than 50% of targets



- Regulations hinge on this question
 - Can we reliably know what's AI generated and what's not?



EU AI Act: first regulation on artificial intelligence



New California Law Will Require AI Transparency and Disclosure Measures



Spain approves registration to regulate AI content and combat deepfakes



- Reuters

OpenAI, Google, others pledge to watermark AI content for safety, White House says



China's Plan to Make Al Watermarks Happen

Watermarking emerges as a promising strategy



- LLMs demonstrate human-level text generation appabilities
 Reuters
- Prone to misuse
 - Misinformation & Fake news
 - Al-generated scams
 - Academic dishonesty
- Regulations hinge on this question
 - Can we reliably know what's Al generated and what's Watern

Watermarking emerged as a promising strategy

OpenAI, Google, others pledge to Watermark AI content for safety, White House says



China's Plan to Make Al Watermarks Happen

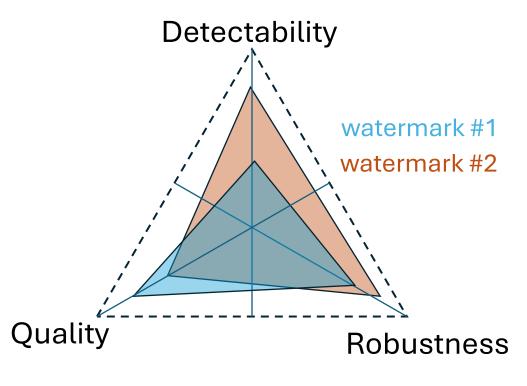
How can we reliably embed & detect watermarks?

Embedding a watermark signal into the token distribution

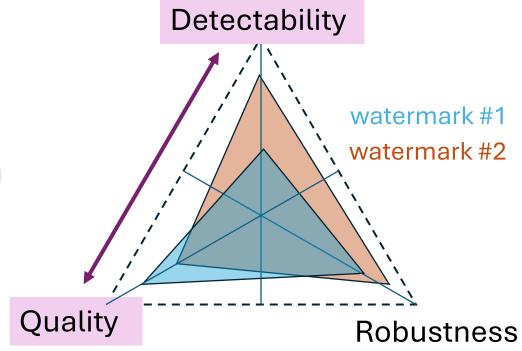
- Embedding a watermark signal into the token distribution
- Myriad watermarking schemes
 - Red-Green (Kirchenbauer et. al. 23)
 - Inverse transform (Kuditipudi et. al. 23)
 - Gumbel-max (Aaronson & Kircher 23)
 - Error correction codes (Christ & Gunn 24)
 - Synth-ID (Dathathri 24)
 - He et. al 24
 - 0 ...

- Embedding a watermark signal into the token distribution
- Myriad watermarking schemes
 - Red-Green (Kirchenbauer et. al. 23)
 - Inverse transform (Kuditipudi et. al. 23)
 - Gumbel-max (Aaronson & Kircher 23)
 - Error correction codes (Christ & Gunn 24)
 - Synth-ID (Dathathri 24)
 - He et. al 24
 - 0

Trade-off between desirable properties



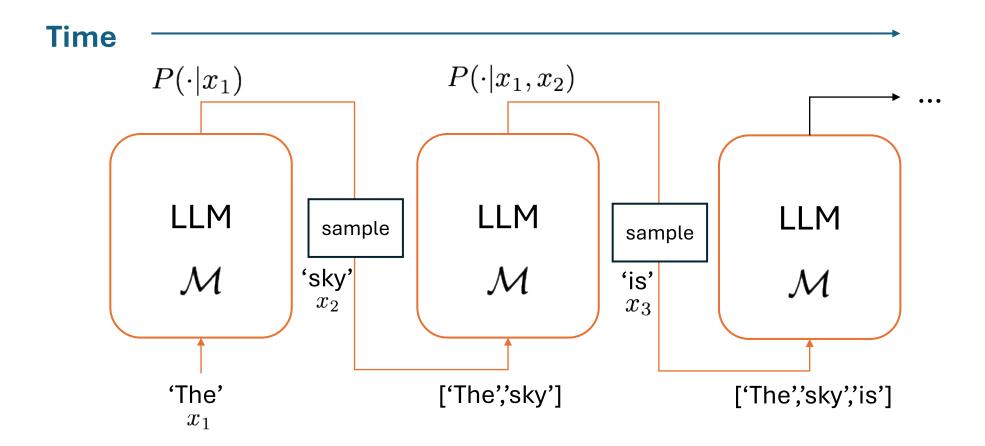
- Embedding a watermark signal into the token distribution
- Myriad watermarking schemes
 - Red-Green (Kirchenbauer et. al. 23)
 - Inverse transform (Kuditipudi et. al. 23)
 - Gumbel-max (Aaronson & Kircher 23)
 - Error correction codes (Christ & Gunn 24)
 - Synth-ID (Dathathri 24)
 - He et. al 24
 - 0



- Trade-off between desirable properties
- We optimize for the quality-detectability trade off.

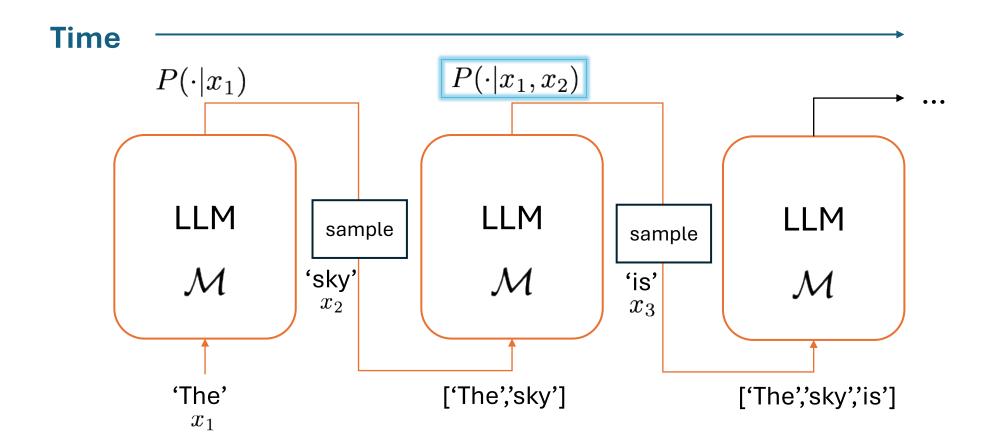
LLM Text Generation

- Tokens, vocabulary $|\mathcal{X}| = m$
- Autoregressive token generation



LLM Text Generation

- We focus on token-level watermarking
 - Treat each distribution independently



- White-box watermarks
 - \circ **Access** to LLM distribution $P(\cdot|x_1,\ldots,x_{i-1})$
 - No access to model weights

- White-box watermarks
 - \circ **Access** to LLM distribution $P(\cdot|x_1,\ldots,x_{i-1})$
 - No access to model weights

Goal: Embed a watermark into the token distribution

- White-box watermarks
 - Access to LLM distribution
 - No access to model weights

Goal: Embed a watermark into the token distribution

Requirements:

- Watermark can be detected from the generated tokens.
- Textual quality is not affected.

- White-box watermarks
 - Access to LLM distribution
 - No access to model weights

Goal: Embed a watermark into the token distribution

Requirements:

- Watermark can be detected from the generated tokens.
- Textual quality is not affected.

Detection – How easy it is to detect the watermark from the text

- White-box watermarks
 - Access to LLM distribution
 - No access to model weights

Goal: Embed a watermark into the token distribution

Requirements:

- Watermark can be detected from the generated tokens.
- Textual quality is not affected.

Detection – How easy it is to detect the watermark from the text

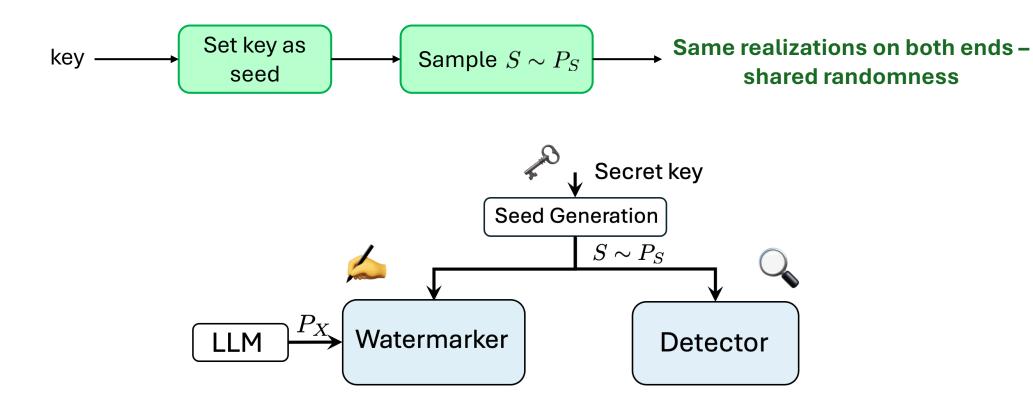
Quality – Distance from original token distribution, distortion

Two parties – Watermarker and Detector

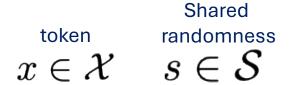


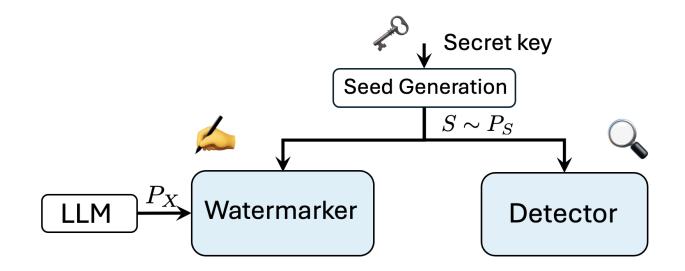
First Step Towards a Formulation

- Both parties share secret key generate Shared randomness
 - \circ Shared side information $S \sim P_S$, we consider uniform over $\mathcal{S} = [0:k-1]$



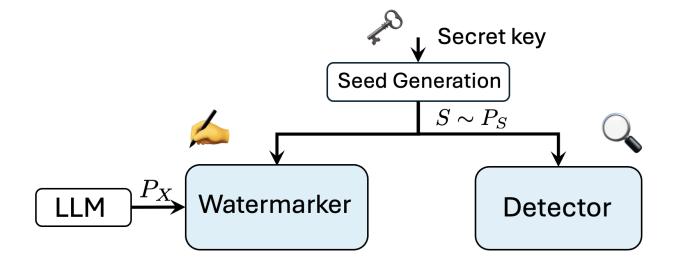
• Score function f(x,s)





Shared randomness

- Score function f(x,s)
 - Watermarker Generate watermarked distribution $P_{X|S}$

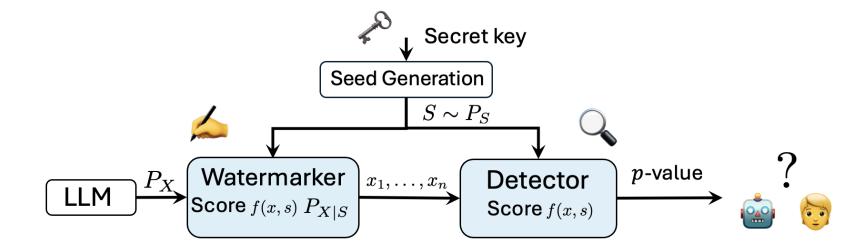


token $x \in \mathcal{X} \quad s \in \mathcal{S}$

token

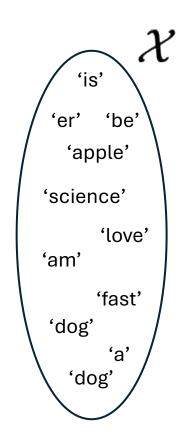
Shared randomness

- Score function f(x,s)
 - Watermarker Generate watermarked distribution $P_{X|S}$
 - **Detector** Apply a statistical test $\frac{1}{n}\sum_{t=1}^{n}f(x_{t},s_{t})\geq \tau$

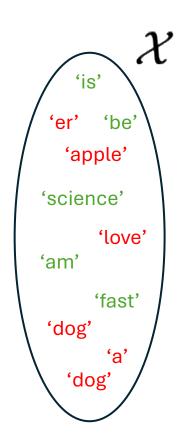


 $x \in \mathcal{X} \quad s \in \mathcal{S}$

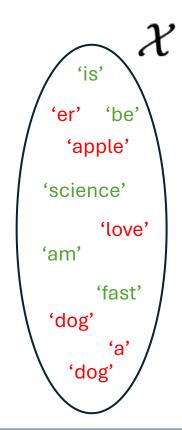
ullet Key idea: Both parties share a partition of token vocabulary ${\mathcal X}$



- ullet Key idea: Both parties share a partition of token vocabulary ${\mathcal X}$
- Binary partition Red list and Green list
 - \circ Random partition according to shared randomness S

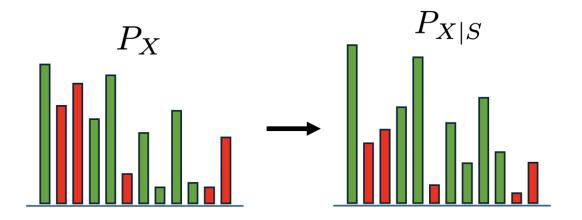


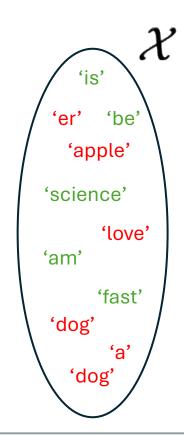
- ullet **Key idea:** Both parties share a partition of ${\mathcal X}$
- Binary partition Red list and Green list
 - \circ Random partition according to shared randomness $oldsymbol{S}$
 - Score function green list membership



$$f(x,s) = \mathbf{1}\{x \in \mathsf{GreenList}(s)\}$$

- **Key idea:** Both parties share a partition of ${\mathcal X}$
- Binary partition Red list and Green list
 - \circ Random partition according to shared randomness S
 - Score function green list membership
- Watermarking Reweigh token distribution
 - Increase the probability of green list tokens
 - Decrease the probability of red list tokens.





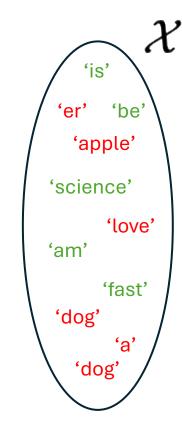
$$f(x,s) = \mathbf{1}\{x \in \mathsf{GreenList}(s)\}$$

$$P_{X|S=s} = \frac{P_X(x)e^{1+\delta f(x,s)}}{\sum_{x'\in\mathbb{X}} P_X(x')e^{1+\delta f(x,s)}}$$

- Statistical test count how many green tokens
 - $_{\circ}$ Detector generates same s samples.

$$\frac{1}{n} \sum_{t=1}^{n} \mathbf{1} \{ x_t \in \mathsf{GreenList}(s_t) \} \ge \tau$$

• Test LHS is a proxy for $\mathbb{E}\left[f(X,S)
ight]$.



$$f(x,s) = \mathbf{1}\{x \in \mathsf{GreenList}(s)\}$$

$$P_{X|S=s} = \frac{P_X(x)e^{1+\delta f(x,s)}}{\sum_{x' \in \mathbb{X}} P_X(x')e^{1+\delta f(x,s)}}$$

Generated tokens are relatively deterministic

```
o Coding P('world' | 'print("hello ')
```

o Math P('3' | '1 + 2 = ')

- Generated tokens are relatively deterministic
 - coding P('world' | 'print("hello ')
 - Math P('3' | '1 + 2 = ')
- Model: min-entropy constraint

$$\mathsf{t} \quad P_{\lambda} = \left\{ P \in \Delta_m : \max_x P(x) \leq \lambda \right\}$$

$$\lambda = \frac{1}{\lambda} \quad \text{Uniform} \qquad \lambda = 1 \quad \text{includes one hot}$$

- Generated tokens are relatively deterministic
 - coding P('world' | 'print("hello ')
 - o Math P('3' | '1 + 2 = ')
- Model: min-entropy constraint

$$\mathsf{t} \quad P_{\lambda} = \left\{ P \in \Delta_m : \max_x P(x) \leq \lambda \right\}$$

$$\lambda = \frac{1}{m} \ \, \Rightarrow \text{Uniform} \qquad \lambda = 1 \ \, \Rightarrow \text{includes one hot}$$

- This work: $\lambda \in \left[\frac{1}{2}, 1\right]$
 - Worst case: A single token gets a probability of 0.5
 - Spiky distributions Token vocabularies are very big.

- Generated tokens are relatively deterministic
 - coding P('world' | 'print("hello ')
 - o Math P('3' | '1 + 2 = ')
- Model: min-entropy constraint

$$\mathsf{t} \quad P_{\lambda} = \left\{ P \in \Delta_m : \max_x P(x) \leq \lambda \right\}$$

$$\lambda = \frac{1}{\lambda} \quad \text{Uniform} \qquad \lambda = 1 \quad \text{includes one hot}$$

- This work: $\lambda \in \left[\frac{1}{2}, 1\right]$
 - Worst case: A single token gets a probability of 0.5
 - Spiky distributions Token vocabularies are very big.

^{*} Are all texts potentially low-entropy?

Optimization Problem Formulation - Objective

Recall:

- Shared randomness s
- $_{\circ}$ Score function f(x,s)
- $_{\circ}$ Watermarked distribution $P_{X|S}$

Optimization Problem Formulation - Objective

Recall:

- $_{\circ}$ Shared randomness s
- $_{\circ}$ Score function f(x,s)
- $_{\circ}$ Watermarked distribution $P_{X|S}$

Watermarked distribution

When X is watermarked: $(X,S) \sim P_S P_{X|S}$

When X is not watermarked: $(X,S) \sim P_S P_X$

Original distribution

Optimization Problem Formulation - Objective

- ullet Recall Proxy for detection Expected score: $\mathbb{E}\left[f(X,S)
 ight]$
- Objective function maximize the gap between expected scores:

$$\mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)]$$

Our contribution: Optimization Problem Formulation

Considerations:

$$\mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)]$$

Our contribution: Optimization Problem Formulation

- Considerations:
 - \circ **Watermarked distribution** chosen after f(x,s) and P_X are set

$$\max_{P_{X,S}} \left[\mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)] \right]$$

Our contribution: Optimization Problem Formulation

- Considerations:
 - $_{\circ}$ Watermarked distribution chosen after f(x,s) and P_{X} are set
 - No control over token distribution

$$\min_{P_X \in P_{\lambda}} \max_{P_{X,S}} \left[\mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)] \right]$$

Our contribution: Optimization Problem Formulation

Considerations:

- Watermarked distribution chosen after score and token are set
- No control over token distribution
- **Score** set prior to 'communication' can't be dependent on anything

$$\max_{f \in \mathcal{F}} \min_{P_X \in P_{\lambda}} \max_{P_{X,S}}$$

$$\max_{f \in \mathcal{F}} \min_{P_X \in P_{\lambda}} \max_{P_{X,S}} \left[\mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)] \right]$$

Optimization Formulation – Zero Distortion

- Zero distortion watermarks: $\mathbb{E}_{P_S}[P_{X|S}] = P_X$
 - Proxy for textual quality
 - Proxy for eavesdropper undetectability

$$\max_{f \in \mathcal{F}} \min_{P_X \in P_{\lambda}} \max_{P_{X,S}} \mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)]$$

Optimization Formulation – Zero Distortion

- Zero distortion watermarks: $\mathbb{E}_{P_S}[P_{X|S}] = P_X$
 - Proxy for textual quality
 - Proxy for eavesdropper undetectability
- Get that 'for free' by optimizing over **couplings** of (P_X, P_S)

$$\left\{ P_{X,S} \middle| \sum_{x \in \mathcal{X}} P_{X,S} = P_S, \sum_{s \in \mathcal{S}} P_{X,S} = P_X \right\}$$

$$\max_{f \in \mathcal{F}} \min_{P_X \in P_{\lambda}} \max_{P_{X,S}} \mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)]$$

Optimization Formulation – Zero Distortion

- Zero distortion watermarks: $\mathbb{E}_{P_S}[P_{X|S}] = P_X$
 - Proxy for textual quality
 - Proxy for textual quality
 Proxy for eavesdropper undetectability

Get that 'for free' by optimizing over **couplings** of
$$(P_X, P_S)$$

$$\left\{P_{X,S}\left|\sum_{x\in\mathcal{X}}P_{X,S}=P_S,\sum_{s\in\mathcal{S}}P_{X,S}=P_X\right.\right\}$$
 Optimal Transport?

$$\max_{f \in \mathcal{F}} \min_{P_X \in P_{\lambda}} \max_{P_{X,S}} \quad \mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)]$$

• Set score function $f \in \mathcal{F}$

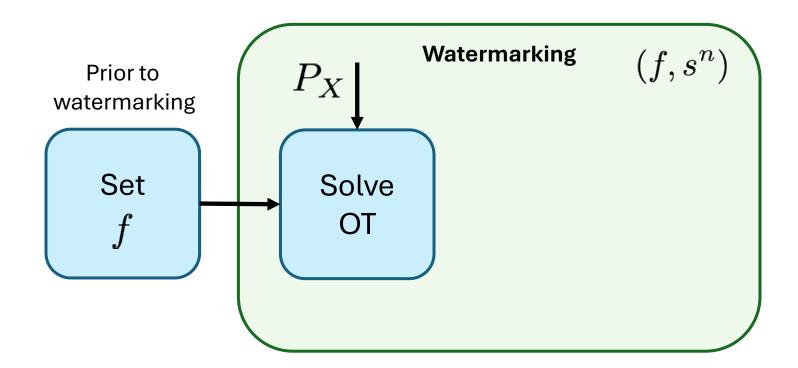
Prior to watermarking

 $egin{pmatrix} \mathsf{Set} \ f \end{pmatrix}$

• Given (f, P_X) solve the Optimal Transport problem:

$$\max_{P_{X,S}} \mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)]$$

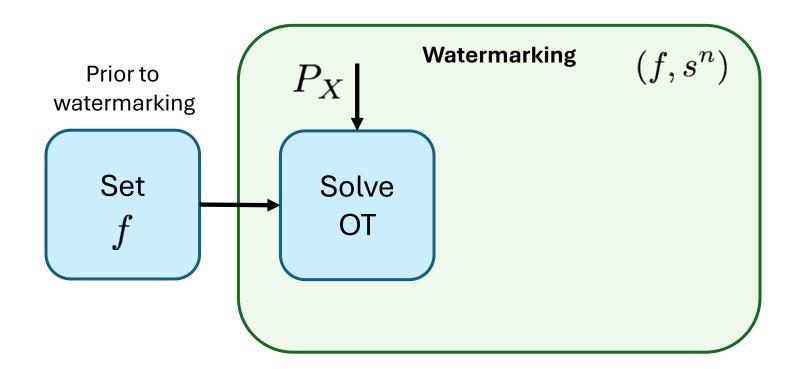
Constant w.r.t coupling



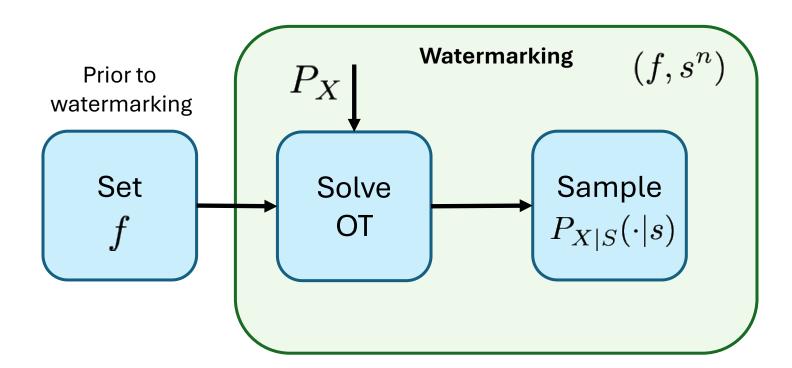
- Given (f, P_X) solve the Optimal Transport problem:
- Sinkhorn

$$\max_{P_{X,S}} \mathbb{E}[f(X,S)]$$

f Is the OT cost

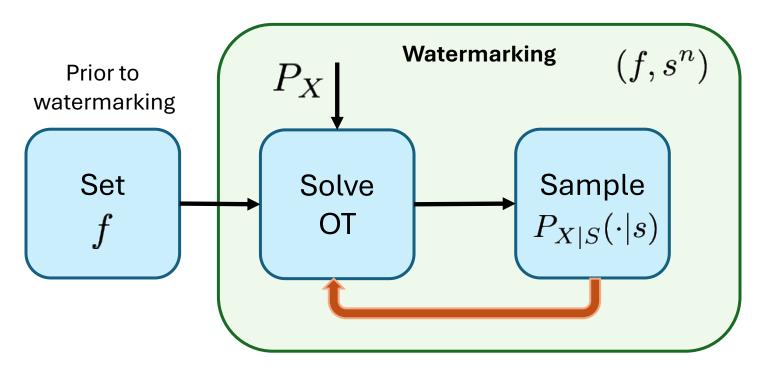


- Extract $P_{X|S}(\cdot|s)$ from Optimal coupling normalized column.
- Sample $P_{X|S}(\cdot|s)$ to obtain a watermarked token



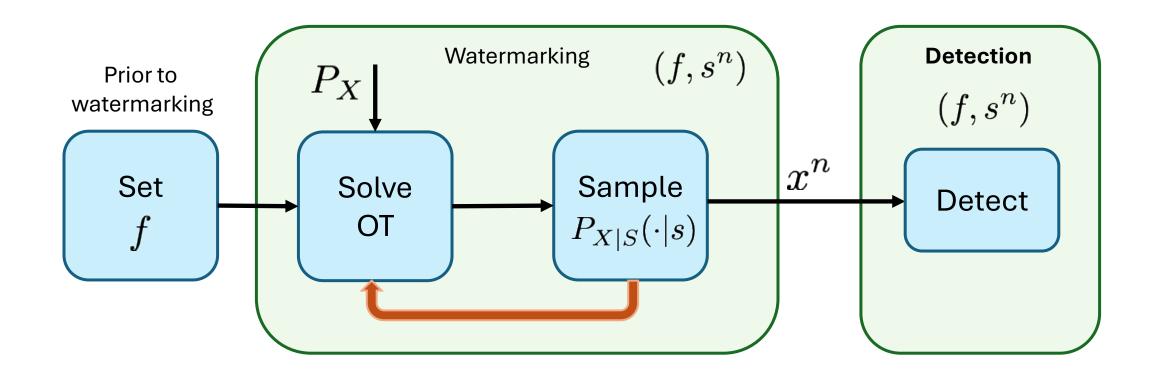
- Extract $P_{X|S}(\cdot|s)$ from Optimal coupling column.
- Sample $P_{X|S}(\cdot|s)$ to obtained watermarked token

Repeat for n tokens



Detector:

 $_{\circ}$ Knows f and performs statistical test: $\dfrac{1}{n}\sum_{t=1}^{n}f(x_{t},s_{t})\geq au$



Optimization Problem Formulation:

$$\max_{f \in \mathcal{F}} \min_{P_X \in P_{\lambda}} \max_{P_{X,S}} \quad \mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)]$$

Optimization Problem Formulation:

$$\max_{f \in \mathcal{F}} \min_{P_X \in P_{\lambda}} \left[\max_{P_{X,S}} \quad \mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)] \right]$$

- Solving an OT via Sinkhorn
 - $_{\circ}$ Score f determines the OT cost
 - Provides a zero-distortion watermark

Optimization Problem Formulation:

$$\mathsf{D}_{\mathsf{gap}}(m,k,\lambda,\mathcal{F}) = \left| \begin{array}{ccc} \max & \min & \max \\ f \in \mathcal{F} & P_X \in P_\lambda & P_{X,S} \end{array} \right. \\ \left. \mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)] \right.$$

- Solving an OT via Sinkhorn
 - $_{\circ}$ Score f determines the OT cost
 - Provides a zero-distortion watermark

• We call the solution the *detection gap* $\mathsf{D}_{\mathsf{gap}}(m,k,\lambda,\mathcal{F})$

Optimization Problem Formulation:

- the solution the **detection gap** $D_{gap}(m, k, \lambda, \mathcal{F})$

SimplexWater Binary Score Functions

```
      0
      0
      0
      0
      0
      0
      0

      0
      0
      0
      1
      1
      1
      1

      0
      1
      1
      0
      0
      1
      1

      0
      1
      1
      1
      0
      0
      1
      0
      1

      1
      0
      1
      1
      0
      1
      0
      1
      0
      1

      1
      1
      0
      1
      0
      1
      0
      1
      0
      1

      1
      1
      0
      1
      0
      0
      1
      0
      1
      0
      1
```

Token vocabulary Shared randomness

• Score class
$$\mathcal{F}_{\mathsf{bin}} = \{f: [1:m] imes [1:k] o \{0,1\}\}$$

- - Example red/green

Token vocabulary Shared randomness

- Score class $\mathcal{F}_{\mathsf{bin}} = \{f: [1:m] imes [1:k] o \{0,1\}\}$
 - Example red/green
- Optimization can be simplified:

Proposition 1. Let $\lambda \in \left[\frac{1}{2},1\right)$. For $f \in \mathcal{F}_{bin}$, define the vector $f_i = [f(i,1),\ldots,f(i,k)] \in \{0,1\}^k$ for each $i \in \mathcal{X}$. Then,

$$\mathsf{D}_{\mathsf{gap}}(m,k,\lambda,\mathcal{F}_{\mathsf{bin}}) = \max_{f \in \mathcal{F}_{\mathsf{bin}}} \quad \min_{i,j \in \mathcal{X}, i
eq j} \quad rac{(1-\lambda)d_H(f_i,f_j)}{k},$$

where $d_H(a,b) = \sum_{i=1}^k \mathbf{1}_{\{a_i \neq b_i\}}$ denotes the Hamming distance between $a, b \in \{0,1\}^k$ and $\mathbf{1}_{\{\cdot\}}$ is the indicator function.

Token vocabulary Shared randomness

- Score class $\mathcal{F}_{\mathsf{bin}} = \{f: [1:m] imes [1:k] o \{0,1\}\}$
 - Example red/green
- Optimization can be simplified:

Proposition 1. Let $\lambda \in \left[\frac{1}{2},1\right)$. For $f \in \mathcal{F}_{\mathsf{bin}}$, define the vector $f_i = [f(i,1),\ldots,f(i,k)] \in \{0,1\}^k$ for each $i \in \mathcal{X}$. Then,

$$\mathsf{D}_{\mathsf{gap}}(m,k,\lambda,\mathcal{F}_{\mathsf{bin}}) = \max_{f \in \mathcal{F}_{\mathsf{bin}}} \quad \min_{i,j \in \mathcal{X}, i
eq j} \quad rac{(1-\lambda)d_H(f_i,f_j)}{k},$$

where $d_H(a,b) = \sum_{i=1}^k \mathbf{1}_{\{a_i \neq b_i\}}$ denotes the Hamming distance between $a, b \in \{0,1\}^k$ and $\mathbf{1}_{\{\cdot\}}$ is the indicator function.

Each element in the vocabulary is assigned with a binary vector

 $\max_{f \in \mathcal{F}} \min_{P_X \in P_{\lambda}} \max_{P_{X,S}} \mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)]$

Binary Scores

- Score class $\mathcal{F}_{\mathsf{bin}} = \{f: [1:m] \times [1:k] \rightarrow \{0,1\}\}$
 - Example red/green
- Optimization can be simplified:

Maximize the *minimum* Hamming distance

Proposition 1. Let $\lambda \in \left[\frac{1}{2},1\right)$. For $f \in \mathcal{F}_{bin}$, define the vector $f_i = [f(i,1),\ldots,f(i,k)] \in \{0,1\}^k$ for each $i \in \mathcal{X}$. Then,

$$\mathsf{D}_{\mathsf{gap}}(m,k,\lambda,\mathcal{F}_{\mathsf{bin}}) = \max_{f \in \mathcal{F}_{\mathsf{bin}}} \quad \min_{i,j \in \mathcal{X}, i
eq j} \quad rac{(1-\lambda)d_H(f_i,f_j)}{k},$$

where $d_H(a,b) = \sum_{i=1}^k \mathbf{1}_{\{a_i \neq b_i\}}$ denotes the Hamming distance between $a,b \in \{0,1\}^k$ and $\mathbf{1}_{\{\cdot\}}$ is the indicator function.

Each element in the vocabulary is assigned with a binary vector

- This is a coding theoretic problem!
- f_i are codewords
- We design a distance-maximizing code

Proposition 1. Let $\lambda \in \left[\frac{1}{2},1\right)$. For $f \in \mathcal{F}_{bin}$, define the vector $f_i = [f(i,1),\ldots,f(i,k)] \in \{0,1\}^k$ for each $i \in \mathcal{X}$. Then,

$$\mathsf{D}_{\mathsf{gap}}(m,k,\lambda,\mathcal{F}_{\mathsf{bin}}) = \max_{f \in \mathcal{F}_{\mathsf{bin}}} \quad \min_{i,j \in \mathcal{X}, i
eq j} \quad rac{(1-\lambda)d_H(f_i,f_j)}{k},$$

where $d_H(a,b) = \sum_{i=1}^k \mathbf{1}_{\{a_i \neq b_i\}}$ denotes the Hamming distance between $a,b \in \{0,1\}^k$ and $\mathbf{1}_{\{\cdot\}}$ is the indicator function.

Maximizing detection gap



(Watermarking problem)

Designing a distance-maximizing code

(Coding theory problem)





Designing a distance-maximizing code

(Watermarking problem)

(Coding theory problem)

import coding_theory

Maximizing detection gap



Designing a distance-maximizing code

(Watermarking problem)

(Coding theory problem)

import coding_theory

Theorem 1 (Maximum Detection Gap Upper Bound). Consider the class of binary score functions \mathcal{F}_{bin} and uniform P_S . Then, for any $\lambda \in \left[\frac{1}{2},1\right)$, the maximum detection gap can be bounded as

$$\mathsf{D}_{\mathsf{gap}}(m,k,\lambda,\mathcal{F}_{\mathsf{bin}}) \leq rac{m(1-\lambda)}{2(m-1)}.$$

 λ - Distribution family parameter

m - Vocabulary size

Maximizing detection gap



Designing a distance-maximizing code

import coding_theory

Theorem 1 (Maximum Detection Gap Upper Bound). Consider the class of binary score functions \mathcal{F}_{bin} and uniform P_S . Then, for any $\lambda \in \left[\frac{1}{2},1\right)$, the maximum detection gap can be bounded as

$$\mathsf{D}_{\mathsf{gap}}(m,k,\lambda,\mathcal{F}_{\mathsf{bin}}) \leq rac{m(1-\lambda)}{2(m-1)}.$$

Proof - Plotkin Bound! (1960)

Achieving the Detection Bound

- Can we achieve the upper bound?
 - Equivalently Design a bound-achieving code

Achieving the Detection Bound

- Can we achieve the upper bound?
 - Equivalently Design a bound-achieving code

Yes! - Simplex codes

Shared randomness size

$$(|\mathcal{S}| = k = m - 1)$$

Definition 1 (Simplex Code). For any $x, s \in [0:m-1]$, let bin(x), bin(s) denote their binary representations respectively using $log_2 m$ bits. A simplex code $f_{sim} : [0:m-1] \times [1:m-1] \rightarrow \{0,1\}$ is characterized by

$$f_{\mathsf{sim}}(x,s) \triangleq \det(\mathsf{bin}(x),\mathsf{bin}(s)),$$

where $dot(bin(x), bin(s)) \triangleq \sum_{i=1}^{\log_2 m} bin(x)_i \cdot bin(s)_i$ and $bin(v)_i$ denotes the ith bit in the binary representation of v.

Achieving the Detection Bound

Simplex code

Definition 1 (Simplex Code). For any $x, s \in [0:m-1]$, let bin(x), bin(s) denote their binary representations respectively using $\log_2 m$ bits. A simplex code $f_{sim}:[0:m-1]\times[1:m-1]\to\{0,1\}$ is characterized by

$$f_{\mathsf{sim}}(x,s) \triangleq \det(\mathsf{bin}(x),\mathsf{bin}(s)),$$
 (5)

where $dot(bin(x), bin(s)) \triangleq \sum_{i=1}^{\log_2 m} bin(x)_i \cdot bin(s)_i$ and $bin(v)_i$ denotes the *i*th bit in the binary representation of v.

 $\bullet \quad \text{Example - Simplex code} \ \ m=4, k=m-1=3$

$$f(1,2) = \mathsf{dot}(\mathsf{bin}(1),\mathsf{bin}(2)) = \mathsf{dot}(001,010) = 0$$

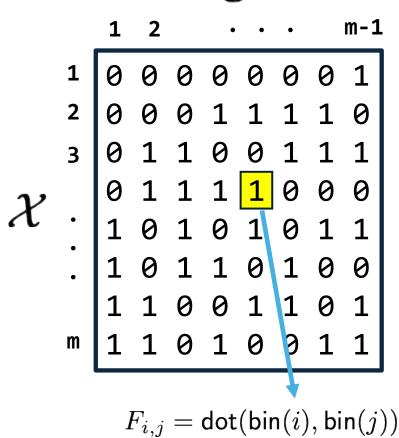
• SimplexWater - Take score f to be a simplex code!

- SimplexWater Take score f to be a simplex code! ${\cal S}$
- $m \times (m-1)$ binary lookup table F

Definition 1 (Simplex Code). For any $x, s \in [0:m-1]$, let bin(x), bin(s) denote their binary representations respectively using $\log_2 m$ bits. A simplex code $f_{sim}:[0:m-1]\times[1:m-1]\to\{0,1\}$ is characterized by

$$f_{\mathsf{sim}}(x,s) \triangleq \det(\mathsf{bin}(x),\mathsf{bin}(s)),$$
 (5)

where $dot(bin(x), bin(s)) \triangleq \sum_{i=1}^{\log_2 m} bin(x)_i \cdot bin(s)_i$ and $bin(v)_i$ denotes the ith bit in the binary representation of v.

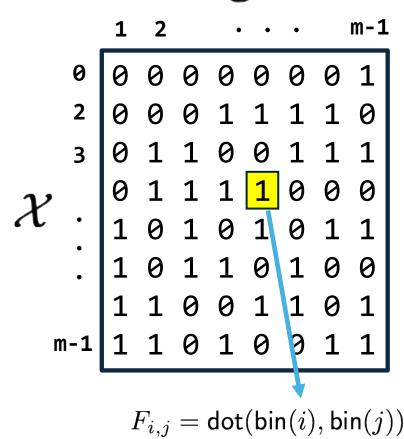


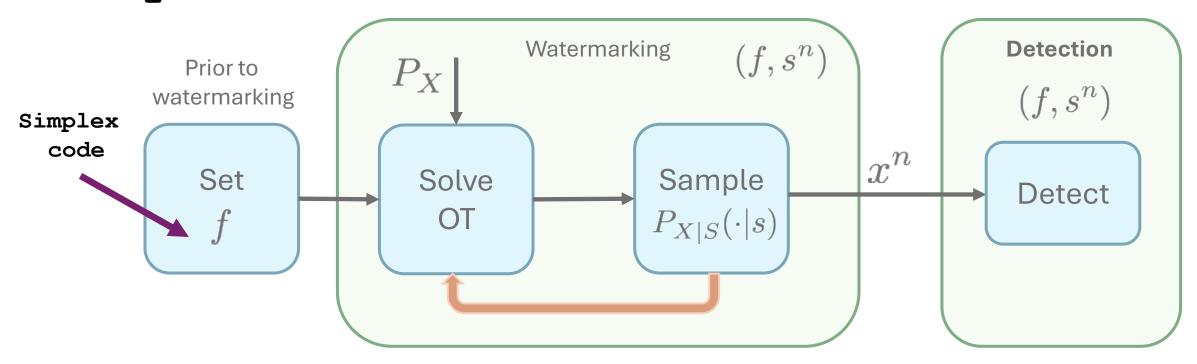
- SimplexWater Take score f to be a simplex code! ${\cal S}$
- $m \times (m-1)$ binary lookup table F
- Computed apriori on both ends

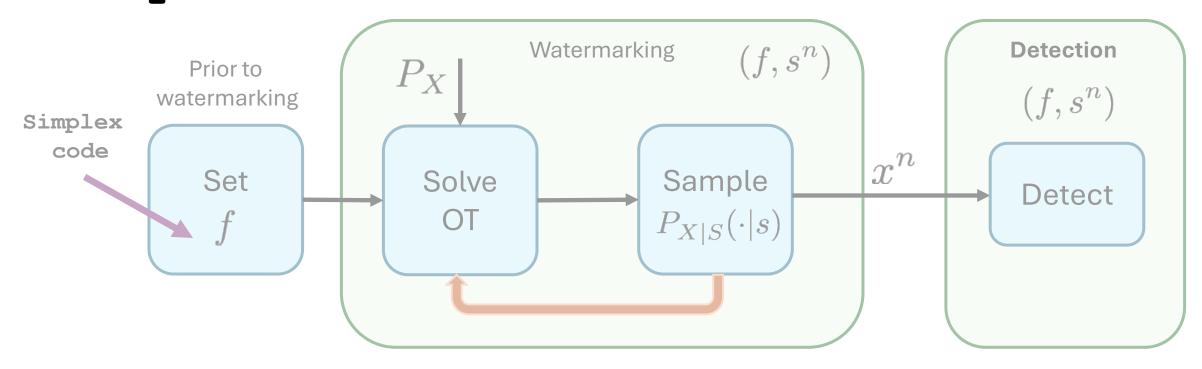
Definition 1 (Simplex Code). For any $x, s \in [0:m-1]$, let bin(x), bin(s) denote their binary representations respectively using $\log_2 m$ bits. A simplex code $f_{sim}:[0:m-1]\times[1:m-1]\to\{0,1\}$ is characterized by

$$f_{\mathsf{sim}}(x,s) \triangleq \det(\mathsf{bin}(x),\mathsf{bin}(s)),$$
 (5)

where $dot(bin(x), bin(s)) \triangleq \sum_{i=1}^{\log_2 m} bin(x)_i \cdot bin(s)_i$ and $bin(v)_i$ denotes the ith bit in the binary representation of v.







• SimplexWater is optimal across all binary score watermarks:

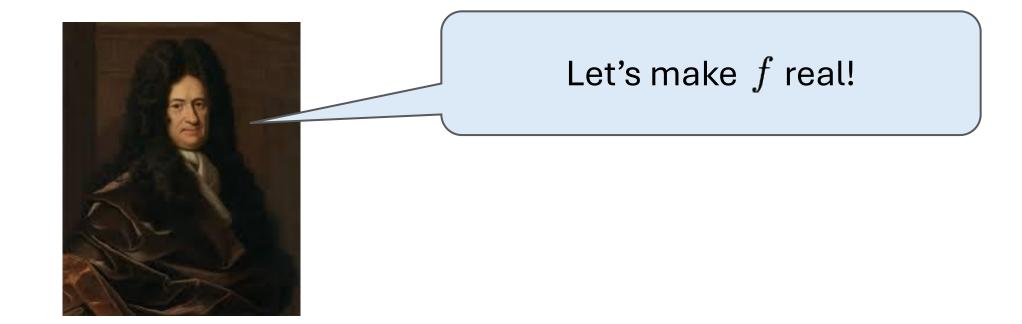
Theorem 2 (SimplexWater Optimality). For any $\lambda \in \left[\frac{1}{2},1\right)$ the maximum detection gap upper bound is attained by SimplexWater.

Beyond Binary Scores

- Binary detection gap is capped by $\ \frac{m(1-\lambda)}{2(m-1)}$
- Bigger field size?

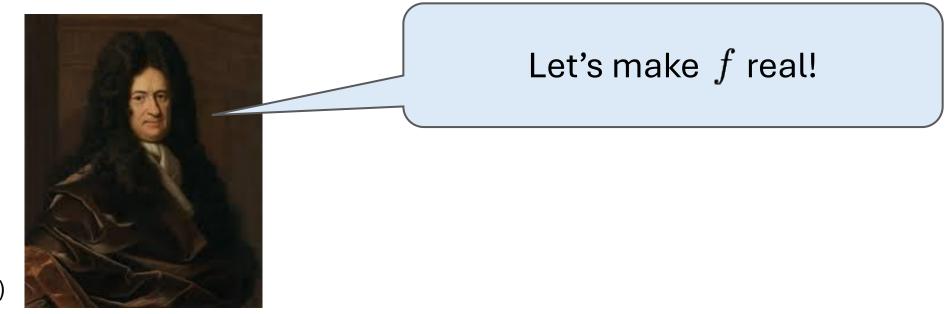
Beyond Binary Scores

- ullet Binary detection gap is capped by $rac{m(1-\lambda)}{2(m-1)}$
- Bigger field size?



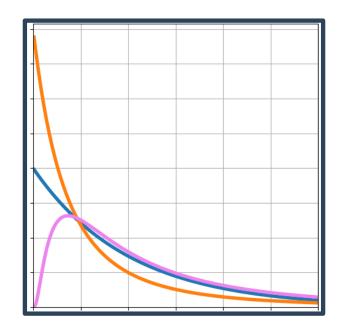
Beyond Binary Scores

- ullet Binary detection gap is capped by $rac{m(1-\lambda)}{2(m-1)}$
- Bigger field size?



*(Leibniz)

HeavyWater Continuous Score Functions

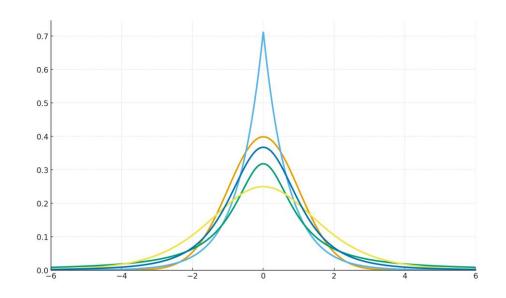


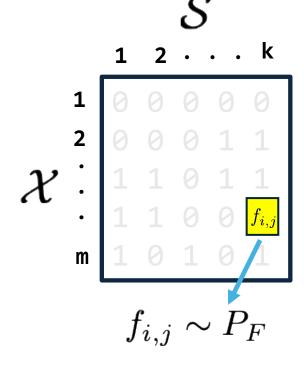
Towards a Continuous Watermark

• Extending continuous scores – Sample f randomly P_F $_\circ$ Practically - Sample the lookup entries i.i.d. P_F

$$f \sim P_F$$

• Distribution P_F is **continuous, zero mean 1D**





Random Score Generalizes Existing Watermarks

- Gumbel Watermark (Aaronson & Kircher '23, OpenAI):
 - Add i.i.d. Gumbel(0,1) noise to the logits and take argmax:

$$x = rg\max_{i \in [1:m]} \ell_i + g_i$$
 (no reweigh)

- o Detection:
 - Resample Gumbel variables and calculate a Gumbel-based statistic



Random Score Generalizes Existing Watermarks

- Gumbel Watermark (Aaronson & Kircher '23, OpenAI):
 - Add i.i.d. Gumbel(0,1) noise to the logits and take argmax:

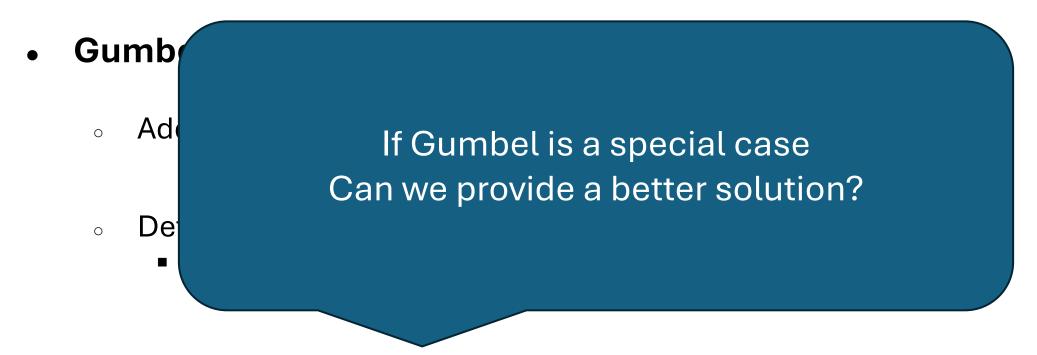
$$x = rg\max_{i \in [1:m]} \ell_i + g_i$$
 (no reweigh)

- Detection:
 - Resample Gumbel variables and calculate a Gumbel-based statistic

Gumbel watermark is asymptotically a special case of our framework!

Theorem 3 (Gumbel Watermark as OT). When the score random variables f(x,s), are sampled i.i.d. from Gumbel(0,1), the solution to the OT problem converges to the Gumbel watermark as $|\mathcal{S}| = k \to \infty$.

Random Score Generalizes Existing Watermarks



Gumbel watermark is asymptotically a special case of our framework!

Theorem 3 (Gumbel Watermark as OT). When the score random variables f(x,s), are sampled i.i.d. from Gumbel(0,1), the solution to the OT problem converges to the Gumbel watermark as $|\mathcal{S}| = k \to \infty$.

What controls the detection gap?

- What controls the detection gap?
- ullet Generally, detection depends on the **tail of** P_F !

Theorem 4 (Detection Gap). Let $\lambda \in \left[\frac{1}{2},1\right)$, and consider the score difference random variable $\Delta = f(x,s) - f(x',s')$ for some $(x,s) \neq (x',s')$, where f(x,s) and f(x',s') are sampled i.i.d. from P_F . Let the cumulative distribution function of Δ be F, and let $Q = F^{-1}$ be its inverse. Then,

$$\lim_{k \to \infty} \mathsf{D}_{\mathsf{gap}}^{[P_F]}(m, k, \lambda) = \int_{1-\lambda}^1 Q(u) du, \tag{8}$$

- What controls the detection gap?
- Generally, detection depends on the **tail of** P_F !

Theorem 4 (Detection Gap). Let $\lambda \in \left[\frac{1}{2},1\right)$, and consider the score difference random variable $\Delta = f(x,s) - f(x',s')$ for some $(x,s) \neq (x',s')$, where f(x,s) and f(x',s') are sampled i.i.d. from P_F . Let the cumulative distribution function of Δ be F, and let $Q = F^{-1}$ be its inverse. Then,

$$\lim_{k \to \infty} \mathsf{D}_{\mathsf{gap}}^{[P_F]}(m, k, \lambda) = \int_{1-\lambda}^1 Q(u) du, \tag{8}$$

- ullet Define score difference Δ
- ullet Quantile of the distribution of Δ is Q
- ullet Detection gap is control by the integral of Q

Theorem 4 (Detection Gap). Let $\lambda \in \left[\frac{1}{2},1\right)$, and consider the score difference random variable $\Delta = f(x,s) - f(x',s')$ for some $(x,s) \neq (x',s')$, where f(x,s) and f(x',s') are sampled i.i.d. from P_F . Let the cumulative distribution function of Δ be F, and let $Q = F^{-1}$ be its inverse. Then,

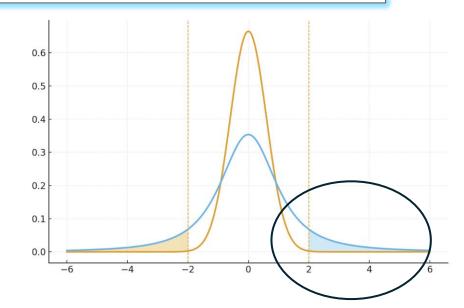
$$\lim_{k \to \infty} \mathsf{D}_{\mathsf{gap}}^{[P_F]}(m, k, \lambda) = \int_{1-\lambda}^1 Q(u) du, \tag{8}$$

• Heavier tail Δ => Bigger detection gap

Theorem 4 (Detection Gap). Let $\lambda \in \left[\frac{1}{2},1\right)$, and consider the score difference random variable $\Delta = f(x,s) - f(x',s')$ for some $(x,s) \neq (x',s')$, where f(x,s) and f(x',s') are sampled i.i.d. from P_F . Let the cumulative distribution function of Δ be F, and let $Q = F^{-1}$ be its inverse. Then,

$$\lim_{k \to \infty} \mathsf{D}_{\mathsf{gap}}^{[P_F]}(m, k, \lambda) = \int_{1-\lambda}^1 Q(u) du, \tag{8}$$

- Heavier tail Δ => Larger detection gap
- Heavy tailed f => Heavy tailed Δ *(through MGFs)

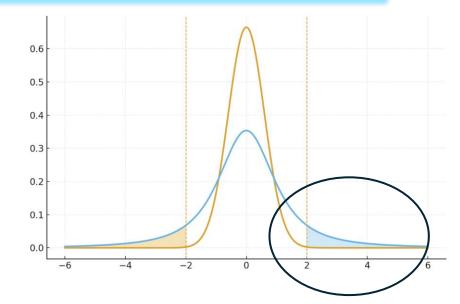


Theorem 4 (Detection Gap). Let $\lambda \in \left[\frac{1}{2},1\right)$, and consider the score difference random variable $\Delta = f(x,s) - f(x',s')$ for some $(x,s) \neq (x',s')$, where f(x,s) and f(x',s') are sampled i.i.d. from P_F . Let the cumulative distribution function of Δ be F, and let $Q = F^{-1}$ be its inverse. Then,

$$\lim_{k \to \infty} \mathsf{D}_{\mathsf{gap}}^{[P_F]}(m, k, \lambda) = \int_{1-\lambda}^1 Q(u) du, \tag{8}$$

Heavier

Heavy t *(throu Sample score from a heavy tailed distribution!

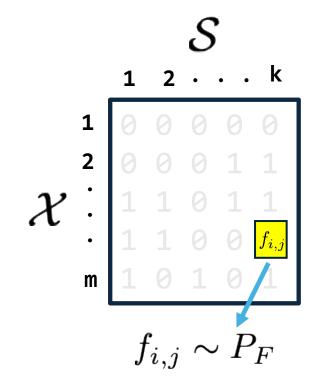


HeavyWater

ullet HeavyWater - Sample f from a heavy tailed distribution P_F

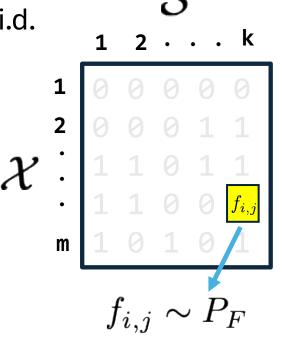
HeavyWater

- ullet HeavyWater Sample f from a heavy tailed distribution P_F
- Best empirical performance Lognormal distribution



HeavyWater

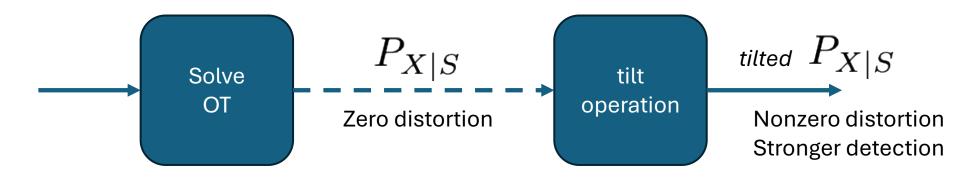
- ullet HeavyWater Sample f from a heavy tailed distribution P_F
- Best empirical performance Lognormal distribution
- HeavyWater Steps:
 - \circ **Set** side information size k and sample a $m \times k$ lookup table i.i.d.
 - Share sampled lookup table prior to text generation
 - Repeat watermarking algorithm



• Recall: $\mathbb{E}_{P_S}[P_{X|S}] = P_X$ (due to OT)

- Recall: $\mathbb{E}_{P_S}[P_{X|S}] = P_X$ (due to OT)
- Higher detectability at the cost of inducing distortion

- Recall: $\mathbb{E}_{P_S}[P_{X|S}] = P_X$ (due to OT)
- Higher detectability at the cost of inducing distortion
- pushing expected scores further apart
- Alter $P_{X|S}$ to **increase** expected score tilting
 - $_{\circ}$ Exact operation $\operatorname{tilt}(P_{X|S})$ depends on exact structure of f



- Tilting SimplexWater:
 - Increase probability of tokens with score 1
 - Decrease probability of tokens with score 0

$$\mathsf{tilt}(P_{X|S=s}^*, s, \delta) = P_{X|S=s}^*(x, s) \left(1 + \delta \cdot (\mathbf{1}_{\{f(x, s)=1\}} - \mathbf{1}_{\{f(x, s)=0\}}) \right)$$

Tilting HeavyWater – around its (zero) mean

$$\mathsf{tilt}(P_{X|S=s}^*,s,\delta) = P_{X|S=s}^*(x,s) \left(1 + \delta \cdot \mathsf{sign}(f(x,s))\right)$$

Numerical Results

Setting

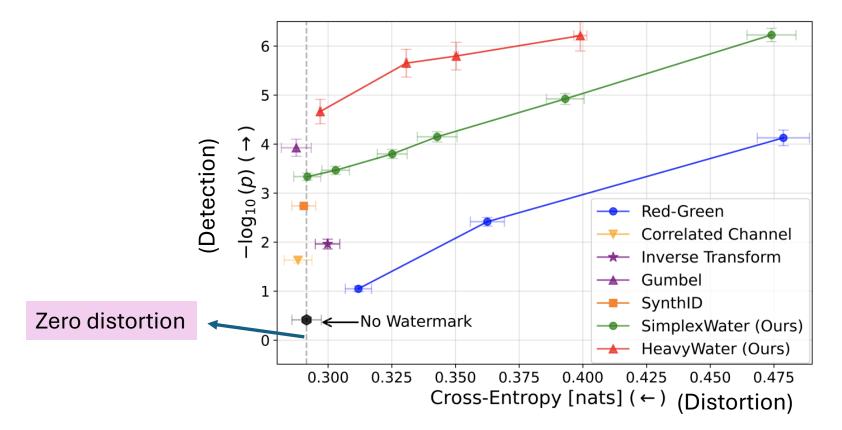
- LLMs Llama2-7b, Llama3-8b, Mistral-v02-7b
- 7 Datasets, For example:
 - FinanceQA
 - LCC
- Top-p sampling, 0.99
- Temperature 1.0/ 0.7
- Methods: Red-Green, Gumbel, Inverse transform, Correlated Channel, SynthID
- Experiments from 2 benchmarks WaterBench and MarkMyWords

Detection-Distortion Tradeoff

- Detection: p-value of statistical test
- Distortion: Cross entropy between base and watermarked distributions

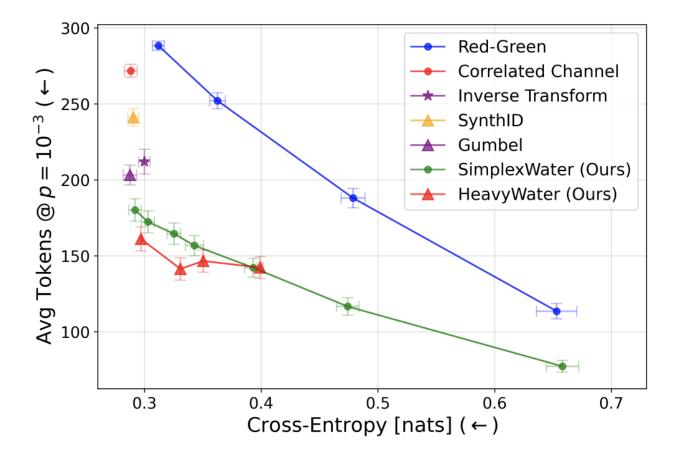
Detection-Distortion Tradeoff

- Detection: p-value of statistical test
- Distortion: Cross entropy between base and watermarked distributions



Watermark Size

Watermark Size - #tokens at detection to obtain a certain p-value.



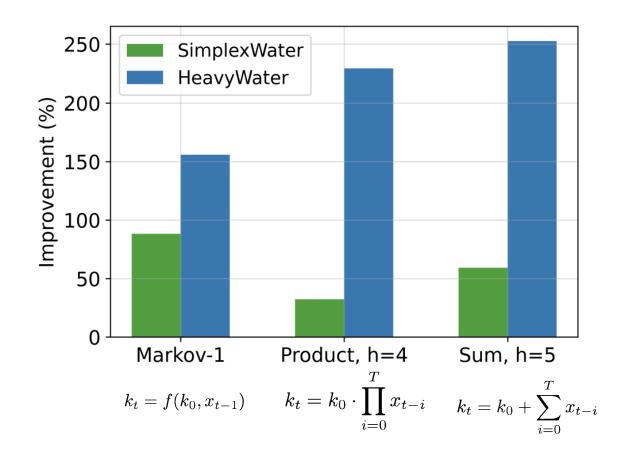
Randomness Generation

- Theoretical assumption Perfect randomness s is i.i.d.
 - New independent key on each time step
- In practice The key is a function of previous tokens
 - Robustness
- Examples of dependence
 - $\mathsf{Markov} \qquad k_t = f(k_0, x_{t-1})$
 - \circ Product $k_t = k_0 \cdot \prod_{i=0}^T x_{t-i}$ \circ Sum $k_t = k_0 + \sum_{i=0}^T x_{t-i}$

• Base method – Red-Green with $\delta=2$ - nonzero distortion

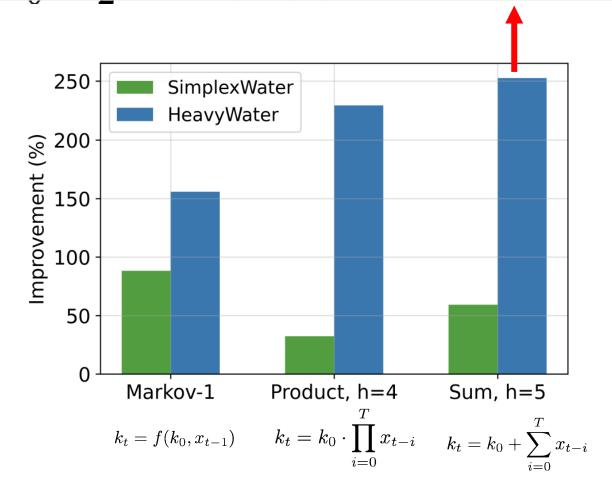
- Base method Red-Green with $\delta=2$ nonzero distortion
- Hashing schemes:
- Markov: $k_t = f(k_0, x_{t-1})$
- Prod: $k_t = k_0 \cdot \prod_{i=0}^T x_{t-i}$
- Sum: $k_t = k_0 + \sum_{i=0}^{T} x_{t-i}$

- Base method Red-Green with $\delta=2$ nonzero distortion
- Hashing schemes:
- Markov: $k_t = f(k_0, x_{t-1})$
- Prod: $k_t = k_0 \cdot \prod_{i=0}^T x_{t-i}$
- Sum: $k_t = k_0 + \sum_{i=0}^{T} x_{t-i}$



Base method – Red-Green wit More than x2 better detection with lower distortion!

- Hashing schemes:
- Markov: $k_t = f(k_0, x_{t-1})$
- Prod: $k_t = k_0 \cdot \prod_{i=0}^T x_{t-i}$
- Sum: $k_t = k_0 + \sum_{i=0}^{T} x_{t-i}$



Watermark Information Size

- How many bits does it take to imprint the watermark?
 - $_{\circ}$ m Tokenizer vocabulary size
 - $_{\circ}$ k side information size (design choice)
 - $_{\circ}$ F floating point precision [bits]

Information efficiency – low resource devices

Watermark Information Size

- How many bits does it take to imprint the watermark?
 - $_{\circ}$ m Tokenizer vocabulary size
 - $_{\circ}$ k side information size (design choice)
 - $_{\circ}$ F floating point precision [bits]

Watermark	# Bits
Red-Green	m
Inverse Transform	$m\mathrm{F}$
Gumbel	$m\mathrm{F}$
SimplexWater (ours)	$\log(m)$
HeavyWater $(ours)$	$\log(k)$

Information efficiency – stealing the watermark (recover key)

Watermark Information Size

- How many bits does it take to imprint the watermark?
 - $_{\circ}$ m Tokenizer vocabulary size
 - $_{\circ}$ k side information size (design choice)
 - $_{\circ}$ F floating point precision [bits]

Watermark	# Bits
Red-Green	m
Inverse Transform	mF
Gumbel	mF
SimplexWater (ours)	$\log(m)$
HeavyWater (ours)	$\log(k)$

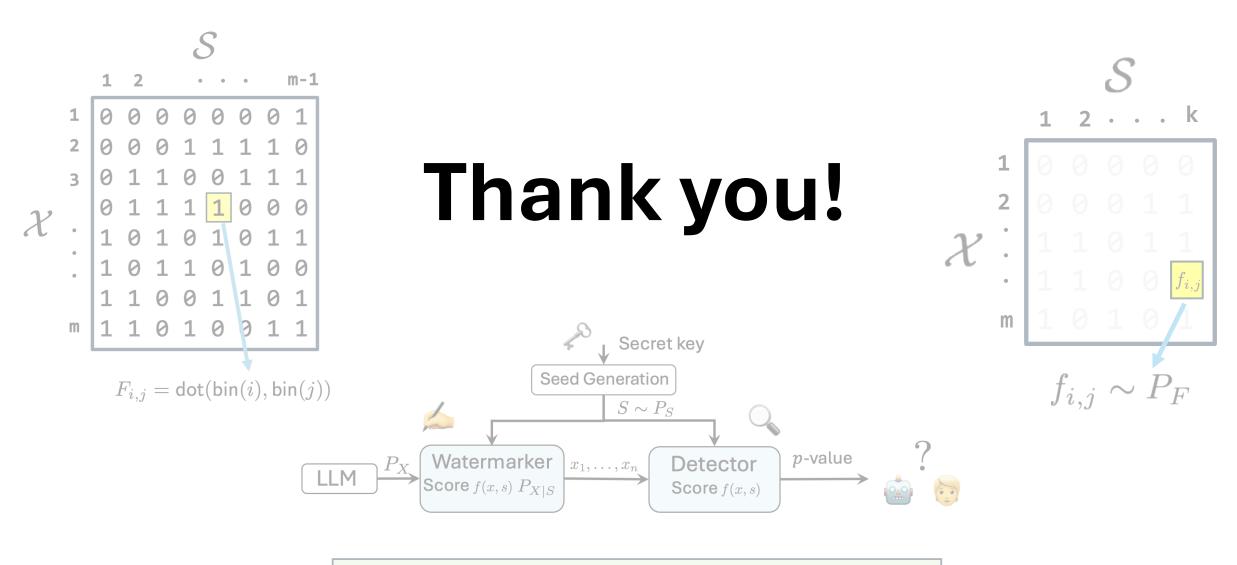
Information efficiency – stealing the watermark (recover key)

Our watermarks are the most information-efficient!

Conclusion

- Main Takeaways

 - Minmax optimization formulation $(f, P_X, P_{X|S})$ Optimal transport solution zero-distortion watermark
 - Score design:
 - Binary scores SimplexWater, Coding theory
 - Continuous scores HeavyWater, Heavy-tailed distributions
 - First unifying framework RG&Gumbel
 - Additional results computational overhead, additional detection&quality metrics, coding.
- Future work
 - Learnable score
 - Integrating robustness into optimization problem



$$\max_{f \in \mathcal{F}} \min_{P_X \in P_{\lambda}} \max_{P_{X,S}} \quad \mathbb{E}_{P_S P_{X|S}}[f(X,S)] - \mathbb{E}_{P_S P_X}[f(X,S)]$$

Randomness Generation

- Theoretical assumption Perfect randomness s is i.i.d.
 - Shared randomness independent across time
- In practice keys are dependent
 - Usually through hashing previous tokens
 - Markov
 - Product Hash
 - Semantic Hash
- We empirically explore the effect of Hashing
- Our watermark works with any randomness generation