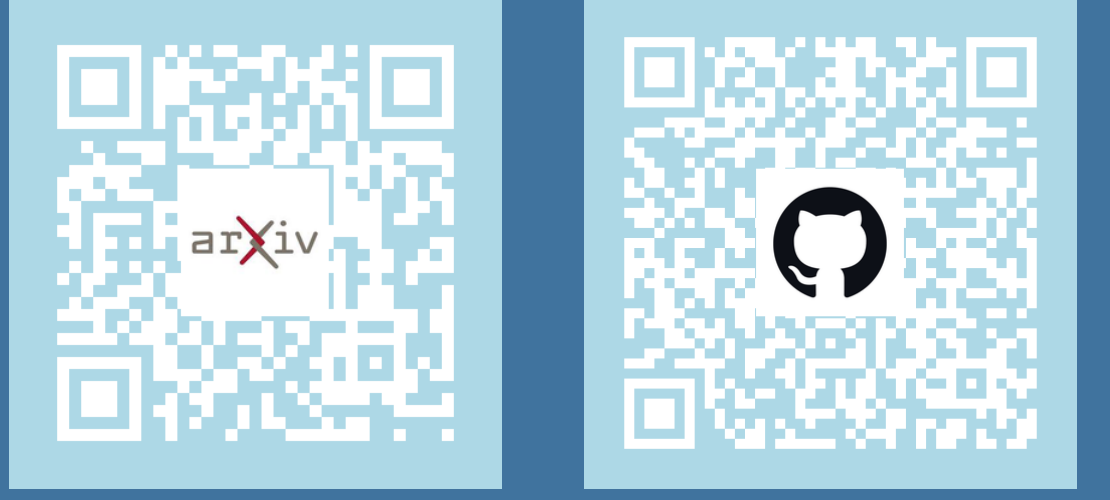


A Statistical Framework of Watermarks for Large Language Models: Pivot, Detection Efficiency and Optimal Rules

Xiang Li¹ Feng Ruan² Huiyuan Wang¹ Qi Long¹ Weijie Su¹

¹ University of Pennsylvania, ² Northwestern University.



Is it possible to (reliably) detect LLM-generated text?

Applications:

- ▶ Fostering original work in education and maintaining academic integrity
- ▶ Preventing fraud and deception
- ▶ Preserving the quality of data for training future AI models

Potential methods: Ad hoc methods leverage context, linguistic patterns, and other markers, which are often inaccurate, unreliable, and biased.

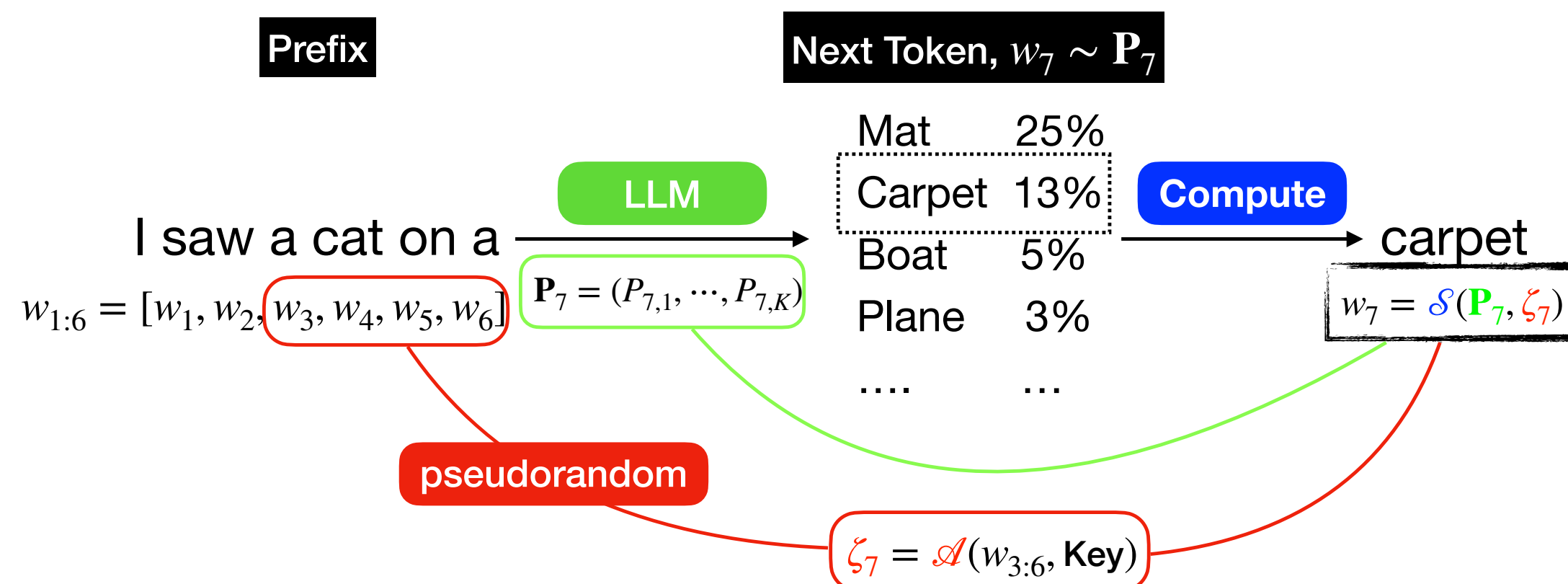
However: Worse, as AI models evolve, LLM-generated text increasingly resembles human-written text!

A principled approach: watermarking LLM-generated text

- ▶ Watermarking embeds subtle statistical signals into LLM-generated text
- ▶ These signal patterns are unlikely to appear in human-written text
- ▶ Watermarking enables provable detection of LLM-generated text

Watermark embedding

- ▶ The vocabulary $\mathcal{W} = \{1, \dots, K\}$, a token w_t , and a text $w_{1:(t-1)} := w_1 \cdots w_{t-1}$.
- ▶ **Autoregressiveness:** An LLM generates each token sequentially by sampling from a probability distribution conditioned on previous tokens:
 $w_t \sim \mathbf{P}_t$ where $\mathbf{P}_t = \text{LLM}(w_{1:(t-1)})$ is next-token prediction (NTP) dist. on \mathcal{W} .
- ▶ **Practical constraint:** \mathbf{P}_t is unknown due to (i) limited access to LLMs and (ii) unknown system/user prompts.

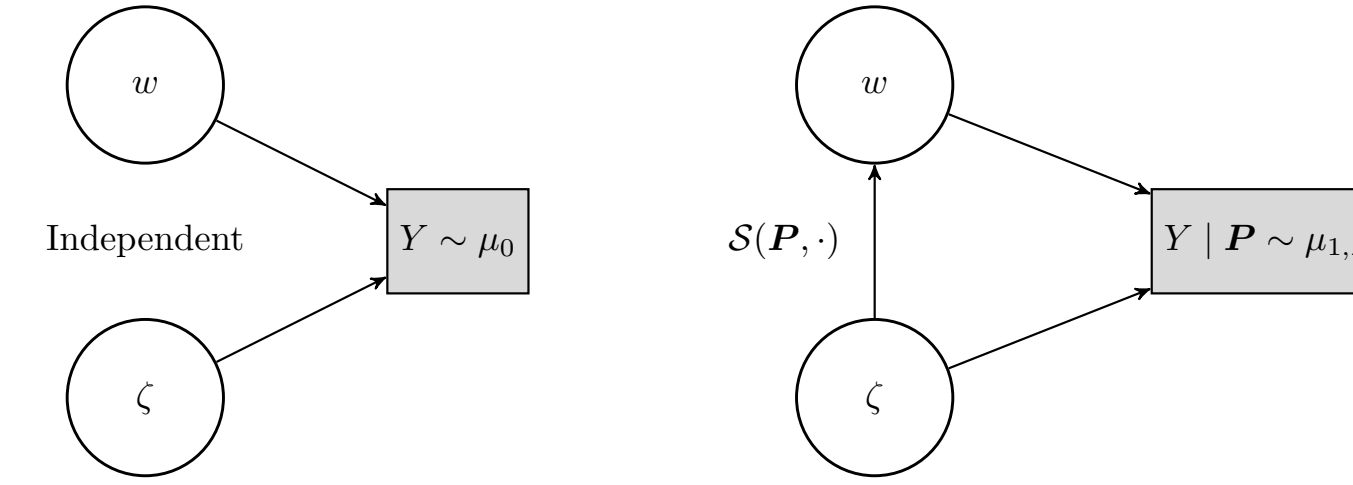


- ▶ Mathematical speaking: $w_t = \mathcal{S}(\mathbf{P}_t, \zeta_t)$ where $\zeta_t = \mathcal{A}(w_{1:(t-1)}, \text{Key})$.
- ▶ $\zeta_{1:t} := \zeta_1 \cdots \zeta_t$ is theoretically i.i.d. and practically recoverable.
- ▶ A watermark is defined by $(\mathcal{A}, \mathcal{S}, \text{Key})$.
- ▶ Watermark signal is the dependence of each w_t on ζ_t .

Watermark detection

Pivotal statistic: Find a scalar function $Y_t = Y(w_t, \zeta_t)$ so that

- ▶ $H_0: Y_t \sim \mu_0$, regardless of \mathbf{P}_t
- ▶ $H_1: Y_t \sim Y(\mathcal{S}(\zeta_t, \mathbf{P}_t), \zeta_t) \stackrel{d}{=} \mu_{1, \mathbf{P}_t}$



Problem formulation: $H_0: Y_t \stackrel{iid}{\sim} \mu_0, \forall t$ vs $H_1: Y_t | \mathbf{P}_t \sim \mu_{1, \mathbf{P}_t}, \forall t$.

Detection rule: Test score $T_h = \sum_{t=1}^n h(Y_t)$ for some score function h . Reject H_0 if T_h is larger than a threshold.

Two considered watermarking schemes

- ▶ A watermark corresponds to sampling from the NTP distribution.
- ▶ **Gumbel-max watermark:** $\zeta = (U_1, U_2, \dots, U_K)$ contains iid $\mathcal{U}(0, 1)$,

$$\mathcal{S}^{\text{gum}}(\mathbf{P}, \zeta) = \arg \max_{w \in \mathcal{W}} \frac{\log U_w}{P_w}, \quad Y^{\text{ars}}(w, \zeta) = U_w.$$

Two scores: $h_{\text{ars}}(y) = -\log(1 - y) > h_{\log}(y) = \log(y)$.

- ▶ **Inverse transform watermark:** Let $\zeta = (\pi, U)$ with $U \sim \mathcal{U}(0, 1)$ and π being sampled uniformly at random from all permutations on \mathcal{W} .

$$\mathcal{S}^{\text{inv}}(\mathbf{P}, \xi) := \pi^{-1}(F^{-1}(U; \pi)), \quad Y^{\text{dif}}(w, \zeta) = |U - \eta(\pi(w))|,$$

where $\eta(w) := \frac{w-1}{|\mathcal{W}|-1}$ and $F(x; \pi) = \sum_{w'} P_{w'} \cdot \mathbf{1}_{\{\pi(w') \leq x\}}$. One uses $h_{\text{neg}}(y) = -y$.

Class-dependent detection efficiency

Questions: (i) How to theoretically compare different score functions and (ii) What is the “optimal” score function?

Class-dependent efficiency: (i) Select a class \mathcal{P} that is believed to contain all $\mathbf{P}_{\leq n}$, and (ii) Evaluate efficiency by the least-favorable power attained over \mathcal{P} .

Prior class: From empirical study, we choose Δ -regular set: $\mathcal{P}_{\Delta} := \{\mathbf{P} = (P_1, \dots, P_K) : \max_w P_w \leq 1 - \Delta\}$.

Formal definition: Fixing Type I error in $(0, 1)$, the pivot-based test statistic $T_h = \sum h(Y_t)$ satisfies

$$\limsup_{n \rightarrow \infty} \sup_{\text{all } \mathbf{P}_t \in \mathcal{P}} [\text{Type II error}]^{\frac{1}{n}} \leq \exp(-R_{\mathcal{P}}(h)),$$

where \mathcal{P} -efficiency rate $R_{\mathcal{P}}(h)$ is defined as

$$R_{\mathcal{P}}(h) = -\inf_{\theta \geq 0} \left\{ \theta \mathbb{E}_0[h(Y)] + \sup_{\mathbf{P} \in \mathcal{P}} \mathbb{E}_{1, \mathbf{P}} \log \left([e^{-\theta h(Y)}] \right) \right\}.$$

Optimal score functions

Finding the optimal score $h^* = \arg \max_h R_{\mathcal{P}}(h)$ reduces to a minimax problem: $\min_h \max_{\mathbf{P} \in \mathcal{P}} L(h, \mathbf{P})$ where $L(h, \mathbf{P}) := \mathbb{E}_0[h(Y)] + \log \left(\mathbb{E}_{1, \mathbf{P}} e^{-h(Y)} \right)$.

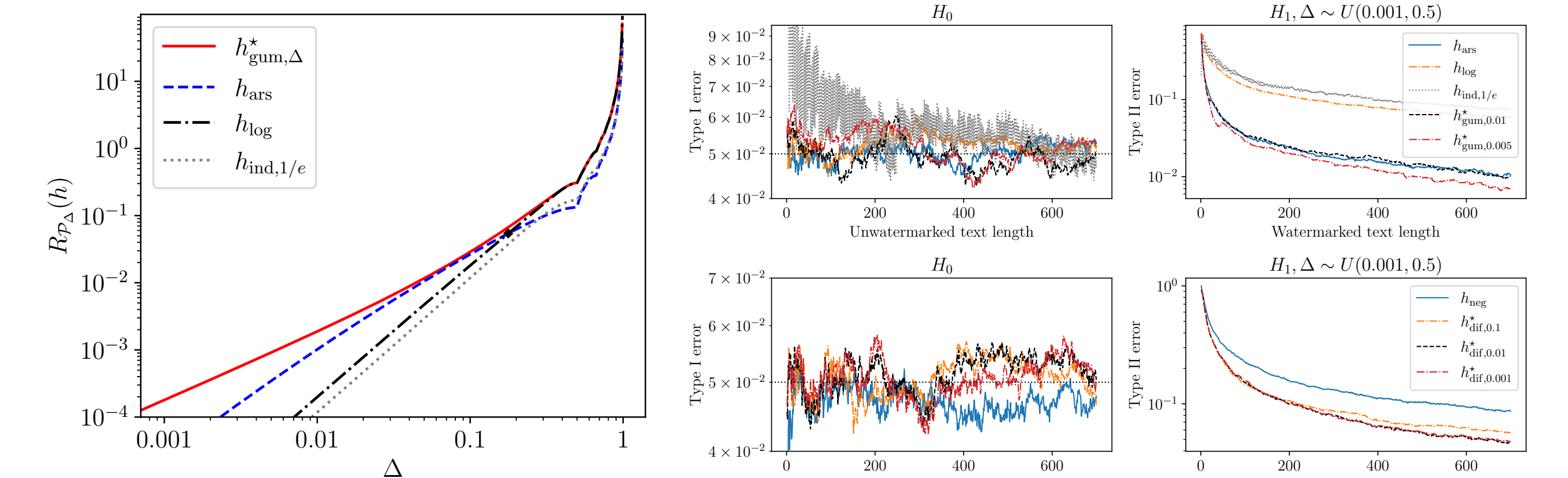
- ▶ The minimax problem is generally not convex-concave. Case-by-case analysis.
- ▶ If there exists an $\mathbf{P}^* \in \mathcal{P}$ and a score function class \mathcal{H} such that for all $h \in \mathcal{H}$,

$$h^* := \log \frac{d\mu_{1, \mathbf{P}^*}}{d\mu_0} \in \mathcal{H}, \quad \sup_{\mathbf{P} \in \mathcal{P}} \mathbb{E}_1[e^{-h(Y)} | \mathbf{P}] = \mathbb{E}_1[e^{-h(Y)} | \mathbf{P}^*],$$

we then have $\max_h R_{\mathcal{P}}(h) = L(h^*, \mathbf{P}^*) = D_{\text{KL}}(\mu_0, \mu_{1, \mathbf{P}^*})$, where the maximum is obtained at h^* .

- ▶ **Main results:** (i) For Gumbel-max, $h_{\text{gum}, \Delta}(y) = \log \frac{d\mu_{1, \mathbf{P}_{\Delta}^*}}{d\mu_0}(y)$ with $\mathbf{P}_{\Delta}^* = (1 - \Delta, \dots, 1 - \Delta, 1 - (1 - \Delta) \cdot \lfloor \frac{1}{1 - \Delta} \rfloor, 0, \dots)$ and (ii) For inverse transform, $h_{\text{dif}, \Delta}^*(r) = \log \frac{f_{\text{dif}, \Delta}(r)}{f_{\text{dif}, 0}(r)}$ where $f_{\text{dif}, \Delta}(r) = \frac{2}{1 - \Delta} \cdot \max \left\{ 1 - \frac{r}{1 - \Delta}, 0 \right\}$ when $K \rightarrow \infty$.

Simulation experiments



LLM experiments

