

Sandcastles in the Storm

Revising the (Im)possibility of Strong Watermarking

Fabrice Harel-Canada* **Boran** Erol* **Connor** Choi Jason Liu Gary Song Violet Peng Amit Sahai





- Watermarking is key to identifying Al content
- "Watermarks in the Sand" [1] argues that every possible watermark can be erased
- We find reason to doubt this theoretical impossibility result... at least for the moment, several watermarking schemes are viable

RQ1 Rapid Mixing दि

Can stationary distributions for watermarking be reached under practical constraints?

IKWYT! Just find the 2nd-largest eigenvalue of the transition matrix, right?

No, the graph of possible responses is $massive \rightarrow computationally intractable$

Universal Attack Formula

• Step 1 (Perturb): A Perturbation Oracle P makes edits (e.g., paraphrases) Step 2 (Check Quality): A Quality Oracle **Q** ensures the edit doesn't degrade quality Step 3 (Repeat): Iterate for sufficiently long to break the watermark. Maybe 200 iterations?





Fact: if mixing occurs, you've reached a stationary distribution + therefore, the "memory" of starting state is *lost*









Every possible response to a prompt is a point in a massive graph

1. P takes a step 2. Q checks if the new state is good enough

Stick to a quality preserving subgraph



Random Walk Attack

watermarked Unwatermarked - topic 1 unwatermarked - topic 2

Semantics can drift so long as the quality stays high!

Key Assumptions

Lineage Distinguisher Tests

Perturbation Oracle	Steps	Tests	Llama-3.1-70B (Failed)	GPT-4o (Failed)	o3-mini-high (Failed)
Word	1000	720	0	0	0
EntropyWord	1000	720	0	0	0
Span	250	720	12	1	0
Sentence	150	720	38	3	0
Document	100	421	2	0	0
Document1Step	100	576	0	0	0
Document2Step	100	678	1	0	0
Total / Failed Tests		4555	53	4	0
Cumulative Distinguished (%)			98.84%	99.91%	100.00%

- Llama3 was a strong and affordable starting point
- Failed tests are sent to the next cheapest model
- Humans are the final boss, but LLMs are good enough



Oracle Reliability

Are LLM-based quality oracles sophisticated enough to guide a random-walk attack?





Takeaways

KA1: Rapid Mixing

Transition probabilities assigned to quality-preserving edits are high



KA2: Reliable Quality Oracle Q is near-perfect to maintain quality throughout the attack



the attack quickly converges to a **stationary** distribution independent of the watermark

too lenient? quality not preserved too conservative? inefficient traversal

Question: Do these assumptions hold up in practice?

Empirical Study Setup

Large-scale empirical study across **718,160** texts **3** watermark schemes, **7** perturbation oracles, **24** quality oracles



Entropy Controlled Prompts

- Vulnerable Domains: Education, Journalism, Creative Writing
- **Progressive Control:** Each prompt more constrained than the last, ex: Lvl 1: "Write a 500-word story" Lvl 2: "...that takes place in Paris"



Watermarkers

- **KGW:** Red/green list based on rolling hash of previous token IDs
- SIR: Uses hash based on semantic embeddings of preceding tokens
- Adaptive: Selectively boosts only high-entropy tokens



OP Precision (%)

RQ3 **Attack Vulnerability**

How effective are random-walk attacks in breaking watermarks when controlling for quality?



Perturbed for *many* steps to ensure sufficient opportunity for mixing

Perturbation Oracles (P)

- **Token:** maskfill random tokens
- **Span:** maskfill contiguous tokens
- **Sentence:** modify a single sentence
- **Document:** full document edits in 1-step, 2-step (modify 1 sentence + global consistency check), multi-step

Quality Oracles (Q)

For original text O and perturbed P:

- **Absolute:** Q scores O / P separately
- **Comparative:** Q sees both O / P together, compares, then scores Many different configurations of oracle

type and LLM base model.

NOTE: Q can be as strong as the watermarking model, but P must be weaker (else just regen with P directly)

[1] Zhang, Hanlin, et al. "Watermarks in the sand: Impossibility of strong watermarking for generative models." ICML (2024).

Rapid mixing is not happening in practice