



The watermarking approach:

- Intentionally embed markers
- Embed using a secret rule (e.g. using a secret key)
- Typically, detection uses the same key.

(although public detection may be possible  
e.g. Fairuze et al. Eprint 2023/1661)

Classical Topic in Cryptography & Security

e.g. watermarking images, text, videos etc.



Very hard. OTOT need to change the text to embed a marker; on the other hand, how do I know that my change hasn't degraded quality?

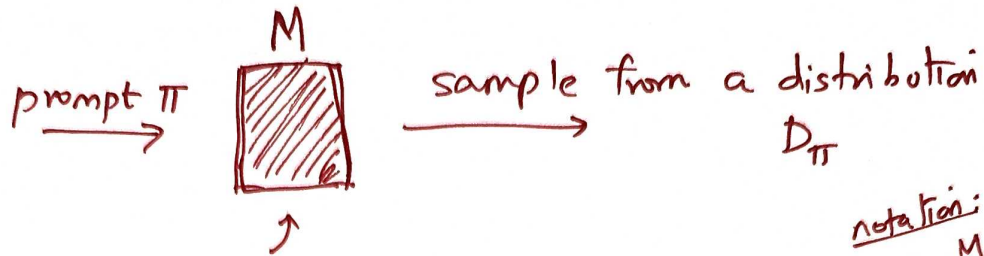
The key Difference for us:

{ We have access not only to an image but the generative process that produced the image.

Gives us Lots of freedom & we will exploit it.

# Generative Models

③



$\uparrow$   
we will treat  
this largely as a black box.

notation:  
Model  $M \Rightarrow D_\pi$

e.g. LLMs (autoregressive),  
Diffusion Models, ...

In fact, today focus on

Autoregressive Models: sample a string from  $T^*$   
 $T$  = set of tokens (for GPT-4  $|T| \sim 10^5$ )

Code for Model  $M$ :

```
t = 0, X_0 = I
while X_t  $\neq$  done
    P_{t+1} = M(\pi, X_1, \dots, X_t)
    sample X_{t+1}  $\sim$  P_{t+1}
    t = t + 1
output X_1, \dots, X_t
```

probability  
distribution  
over  $T$

neural network  
model

e.g.  $\pi(X_1, \dots, X_6) =$  "As a language model, I"

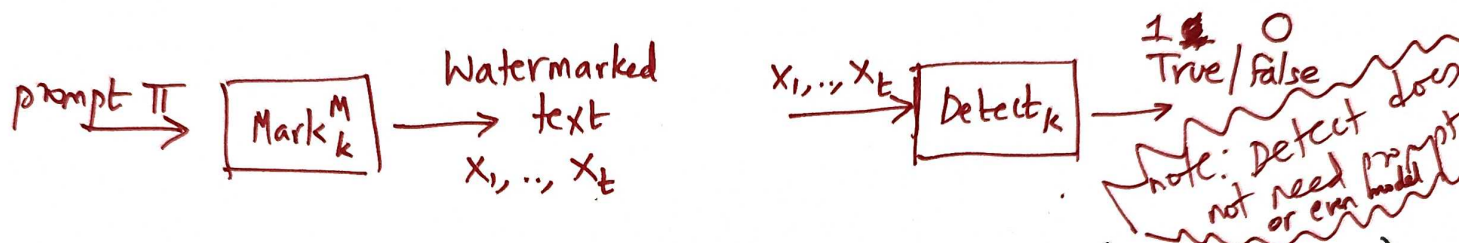
$P_7 =$  "can": 10%.  
"cannot": 70%.  
"am": 5%.  
"do": 3%.  
...

# Watermarking LLMs (= autoregressive models)

④

The "API": Two algorithms  $\text{Mark}_k^M$ ,  $\text{Detect}_k$ :

- Mark and Detect share a private key  $k$
- Mark knows the (autoregressive) model  $M$



TBD: What properties do we want from (Mark, Detect)?

## Strawman 1:

Recall Pseudorandom functions  $\mathcal{F} = \{F_k: \{0,1\}^n \rightarrow \{0,1\}^m\}$   
s.t. for all p.p.t. algorithms  $A$ ,

Exercise:

$$\text{key} \geq \begin{cases} \lambda \text{ bits} \\ \text{long} \end{cases} \left\{ \left| \Pr[A^{F_k} \rightarrow 1] - \Pr[A^R \rightarrow 1] \right| \leq \text{negl}(\lambda) \right\} 2^{-\lambda}$$

where  $\lambda$ : security parameter  
 $R$ : uniformly random function

PRFs give us a good (comp. secure) secret-key encryption scheme:

$$\left\{ \begin{array}{l} \text{Enc}_k(m; r) = (r, \text{PRF}_k(r) \oplus m) \\ \text{Dec}_k(r, c) = c \oplus \text{PRF}_k(r) \end{array} \right.$$

PRFs also give us a good <sup>message</sup> authentication scheme (MAC). (5)

$$\text{MAC}_k(m) = \text{PRF}_k(m)$$

$$\text{Verify}(m, t) = (\text{YES iff } t = \text{PRF}_k(m))$$

Security: given  $(m_1, t_1), \dots, (m_k, t_k)$ , computationally (unforgeability) hard for any p.p.t. algorithm to output  $(m^*, t^*)$  where (a)  $m^* \notin \{m_1, \dots, m_k\}$  AND (b)  $t^* = \text{MAC}_k(m^*)$

$\text{Wat}_k^M(\pi)$  = generate  $x \sim M(\pi)$   
output  $(x, \text{MAC}_k(x))$

- Watermarked texts are correctly recognized (completeness)
- For any given text  $x$  and tag  $t$ , the chance that  $t = \text{MAC}_k(x)$  "accidentally" (i.e. for a random  $k$ ) is small ( $\leq 2^{-x}$ ).  
(soundness / no false positives)

✗ trivially removable.

Lesson: Watermarks cannot be trivially appended. Have to be intrinsic.

## strawman 2: the red/green scheme

(6)

- Randomly partition the tokens into RED and GREEN tokens.  
(Apple, Alphabet, Arugula, Kale, Bacon, Cheese)
- Increase the probability of the words in the green list.  
(in an extreme, only output green tokens)
- Detection is easy (e.g. count # green tokens)

### Problems?

- 1) Quality Degrades: Model prefers not to talk about ~~bacon~~ apples.
- 2) If you talk about bacon or cheese a lot, you may be falsely accused: no soundness.
- 3) When is the watermark detectable?

completeness,

assuming entropy

if output  $\pi$  ~~is~~ is deterministic,

(the quick brown fox .....

no hope to watermark: either change the output

& lose quality or

don't change the output

& lose soundness.

Property 1: Soundness

(no false positives)

— Important! absolute!

(7)

for all text  $t$ ,  $\Pr_{sk} [\text{Detect}_{sk}(t) = \text{TRUE}] \leq 2^{-\lambda}$ .

we won't think about this too much

Stronger: unforgeability

Given watermarked  $t_1, t_2, \dots, t_k$  cannot produce  $t^* \notin \{t_1, \dots, t_k\}$  s.t.  $\text{Detect}_{sk}(t^*) = \text{TRUE}$

Property 2: Completeness

(b-completeness)

— conditional on enough entropy

Text generated by watermarked model will be recognized (given  $sk$ ).

for all prompts  $\pi$

$\Pr_{sk} \left[ \begin{array}{l} \text{Detect}_{sk}(t) = \text{False} \text{ AND} \\ \text{Entropy}(\pi) \geq b \end{array} \right] \leq 2^{-\lambda}$   
we will need to refine

stronger: robustness (w.r.t. modifications  $\mathcal{E}$ )

..... even after modified using transformation  $\mathcal{E}$

for all prompts  $\pi$

$\Pr_{sk} \left[ \begin{array}{l} \text{Detect}_{sk}(\mathcal{E}(t)) = \text{False} \text{ AND} \\ \text{Entropy}(\pi) \geq b \end{array} \right] \leq 2^{-\lambda}$

A very difficult property to achieve in general

we will in next lecture look at substitutions, edits...

Property 3: "Quality"

Very hard to define but what if output of watermarked model is "as good as" the unwatermarked model for all <sup>(poly time)</sup> applications?

↳ Undetectability

For all models & for all efficient algorithms A,

$$\left| \Pr [A^{\text{Model}^M} \rightarrow 1] - \Pr [A^{\text{Mark}_{sk}^M} \rightarrow 1] \right| \leq \epsilon.$$

weaker notion: "distortion-freeness"

For all prompts  $\pi$  and models M,

$$\text{Mark}_{sk}^M \equiv \text{Model}^M(\pi)$$

"Watermarked dist. same as original distribution for a single sample".

... next page

# An undetectable watermarking scheme (Christ-Gunn-Zamir 2024)

- For simplicity & w.l.o.g. think of token space =  $\{0, 1\}$
- so, given a ~~prefix~~ prompt  $\pi$  & prefix  $x_1, \dots, x_{t-1}$ , the model defines a probability  $P_t$  s.t

$$x_t \sim \text{Binomial}(P_t)$$

$$(\Pr(x_t=1) = P_t)$$

## Scheme:

IDEA  
Correlated sampling

key  $K = (u_1, u_2, \dots, u_T) \in [0, 1]^T$  ( $T$  random real numbers in  $[0, 1]$ )

Mark  $K^M$ : for each  $t$ , output  $\begin{cases} 1 & \text{if } u_t \leq P_t \\ 0 & \text{otherwise} \end{cases}$

- Marginal dist of  $x_t$  is Binomial( $P_t$ ) indeed.
  - However,  $x_t$  and  $u_t$  are correlated. (for a fixed  $P_t$ , the smaller  $u_t$ , the more likely  $x_t$  is 1)
  - OTOH, if  $P_t=0$  or  $P_t=1$ ,  $u_t$  has no effect on  $x_t$ . (no entropy)
- PERFECT QUALITY
- COMPLETENESS & SOUNDNESS

Detect $_{K^M}(x_1, \dots, x_T)$ : computes a score  $s(x_j, u_j) = \begin{cases} -\ln u_j & \text{if } x_j=1 \\ -\ln(1-u_j) & \text{if } x_j=0 \end{cases}$

$$c(x) = \sum_{j=1}^T s(x_j, u_j)$$

if  $c(x)$  is large, output watermarked.  
else output non-watermarked.

QUICK CALCULATION

Completeness & Soundness: for any fixed text  $x_1, \dots, x_T$  (not correlated with  $u_1, \dots, u_T$ )

each  $s(x_j, u_j)$  is an exponential random var.

$$\left( \begin{array}{l} -\log U_{[0,1]} \\ \text{pdf}(x) = e^{-x} \quad x \geq 0 \end{array} \right)$$

$$\mathbb{E}_{u_j} [s(x_j, u_j)] = \int_0^1 \ln(1/u) du = 1$$

$$\Rightarrow \mathbb{E}_{u_1, \dots, u_T} [c(x) - \underbrace{|x|}_{=T}] = 0$$



for watermarked outputs,

~~$\mathbb{E}_{u_j} [s(x_j, u_j) | x_j = 1]$~~   
 ~~$\Pr(x_j \neq 1)$~~

~~$\mathbb{E}_{u_j} [s(x_j, u_j)] = \int_0^1 \ln(1/u) du$~~

$$\int_0^1 s(x(u), u) du = \int_0^{P_j} \mathbb{1}[u \leq P_j] \cdot \ln(1/u) + \mathbb{1}[u > P_j] \cdot \ln(1/(1-u)) du$$

$$= \int_0^{P_j} \ln(1/u) du + \int_{P_j}^1 \ln(1/(1-u)) du$$

$$= \int_0^{P_j} \ln(1/u) du + \int_0^{1-P_j} \ln(1/u) du$$

(11)

$$= p_j - p_j \ln p_j + (1-p_j) - (1-p_j) \ln(1-p_j) \quad \left| \int \ln u = u \ln u - u \right.$$

$$= 1 + \ln(2) \cdot H(p_j)$$

**Detector**  
 if  $c(x) - |x| \geq \lambda \sqrt{|x|}$   
 say "watermarked"  
 else "not watermarked"

So,

$$\mathbb{E}_{u_1, \dots, u_T} [c(x) - |x|] = \ln 2 \cdot H(\underbrace{\text{Model}(\pi)}_{\text{distribution on input prompt } \pi})$$

**Problem 1** Expectations different but doesn't guarantee Completeness & soundness

Need to analyze e.g. variance.

(what if high shannon entropy but outputs a fixed  $(x_1, \dots, x_T)$  with prob  $1/2$ ?)

(is the  $1/2$  prob output watermarked?   
 = output by  $\text{Mark}_k^M$  & detected by  $\text{Detect}_k^M$ ?)  
 if so, no soundness.

if not, no ~~completeness~~ Undetectability or no completeness if model decides not to output it

Way out: require completeness only for outputs with high empirical entropy

$$\left| \begin{aligned} \text{Emp. } H^{M, \pi}(x) &= -\log \Pr[\text{Model}^M(\pi) = x] \\ H(\pi) &= \mathbb{E}_{x \sim \text{Model}(\pi)} [\text{Emp. } H^{M, \pi}(x)] \end{aligned} \right.$$

Problem 2 Long keys

Use pseudorandomness  
key = k      $u_j = \text{prf}_k(j)$

once you fix  $(u_1, \dots, u_T)$   
model is deterministic!

Problem 3 multiple outputs

Can't use the same  $u_j \rightarrow$  not undetectable

how about  $u_j = \text{prf}_k(\pi || j) \rightarrow$  detector doesn't know  $\pi$   
& what if  $\pi$  repeats?

how about pick random nonce  $r$   
and let  $u_j = \text{prf}_k(r || j)$

Concatenate  $r$  to model output  $\rightarrow$  same issue as MAC

Solution: generate the first  $l$  tokens until  
they have accumulated enough entropy  
i.e. until  $\sum_{i=1}^l -\log \Pr(x_i | x_1, \dots, x_{i-1}) \geq \lambda$ .

- use  $x_1, \dots, x_l$  as "seed."      $-\log \Pr(x_1, \dots, x_l)$
- set  $u_j = \text{prf}_k(x_1, \dots, x_l || j)$

Detector does not know  $l$  so it tries all possibilities ( $l=1, 2, 3, \dots, T$ )

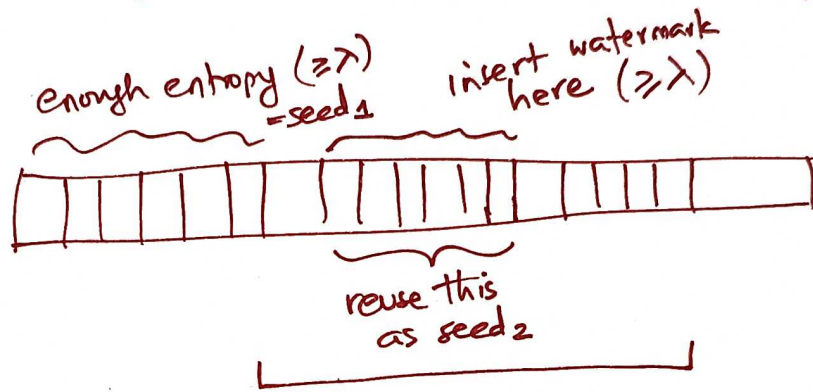
problem 3 unremovability?

:( easily removed if <sup>many</sup> ~~one~~ of the seed bits is changed

a fix: substring-robustness:

verification only needs a sufficiently long contiguous substring

Idea:



as long as some long enough contiguous stry is observed, watermark is detectable.

Typo's from here on !!! → Refer to CGZ for formal proofs!

Additional Material: full construction + proof sketch (without substring-robustness)

Mark<sub>k</sub><sup>M</sup>(π)

i=1  
H=0

while done ≠ (x<sub>1</sub>, ..., x<sub>i-1</sub>) do:

P<sub>i</sub> = Model(π, x<sub>1</sub>, ..., x<sub>i-1</sub>);

if H < λ // if not enough entropy yet

sample x<sub>i</sub> ← P<sub>i</sub> // collect more entropy

H = H + (-log p<sub>i</sub> · x<sub>i</sub> + (1-p<sub>i</sub>)(1-x<sub>i}))  
// how much add'l entropy did I accumulate</sub>

if H ≥ λ then  
r = (x<sub>1</sub>, ..., x<sub>i</sub>) // r is the input to the PRF

else

x<sub>i</sub> ← 1 if F<sub>sk</sub>(r, i) ≤ p<sub>i</sub>  
0 o/w // embed watermark

end.

i = i + 1

end.

so enumerate over all choices  
Detect (x<sub>1</sub>, ..., x<sub>T</sub>):  
don't know how long PRF input is  
for i ∈ [T] do

r<sup>(i)</sup> = (x<sub>1</sub>, ..., x<sub>i</sub>)

Define  
V<sub>j</sub><sup>(i)</sup> = F<sub>sk</sub>(r<sup>(i)</sup>, j)  
if x<sub>j</sub> = 1  
1 - F<sub>sk</sub>(r<sup>(i)</sup>, j)  
if x<sub>j</sub> = 0

if  
∑<sub>j=i+1</sub> ln(1/V<sub>j</sub><sup>(i)</sup>)  
> (L-i) + λ√(L-i)  
return true  
~~end.~~  
end.

if score is high enough, decide "text is watermarked!"

For analysis, assume  $F_{sk}$  is a random function, invoking the pseudorandomness guarantee + fact that only black box use of  $F_{sk}$  + violation of completeness/soundness/undetectability is testable.

Soundness Fix some text  $x_1, \dots, x_T$  (ind of  $sk$ ).

Want:  $\Pr [\text{Detect}(x) = \text{TRUE}] \leq 2^{-\lambda}$ .

For  $i, j \in [T]$  define (1)  $r_j^{(i)}$  (2)  $u_j^{(i)} = O(r_j^{(i)}, j)$   $\leftarrow O$  is a random fn.  
$$v_j^{(i)} = \begin{cases} u_j^{(i)} & \text{if } x_j = 1 \\ 1 - u_j^{(i)} & \text{if } x_j = 0 \end{cases}$$

$u_j^{(i)}$  is independent of  $x_j$  so

$v_j^{(i)} \sim U([0, 1])$

$E_j^{(i)} = \ln(1/v_j^{(i)})$  are exponential random variables w/ rate 1.

Fix  $i$ :  $\Pr \left[ \sum_{j=i+1}^L E_j^{(i)} > (L-i) + \lambda \sqrt{L-i} \right] \leq \dots$

Union bound over  $i \in [T]$

$\Pr \left[ \exists i \in [T] \text{ st } \dots \right] \leq T \cdot 2^{-O(\lambda)} = \text{negligible in } \lambda$

# Undetectability

~~Want~~ (III)  
 Given a fixed history of  $(\pi_i, x_i)$  pairs,  $(i < l)$ , we will show that  $\text{Wat}_M^0(\pi)$  is stat. indistinguishable from  $\text{Model}^M(\pi)$ .

if  $r_i = \text{prefix}(x_i)$  is distinct from

$$r = \text{prefix}(\text{Model}^M(\pi))$$

then  $O(r, *)$  is an independent random function from

$$O(r_i, *) \text{ for all } i$$

$$\text{Therefore } \text{Wat}_M^0(\pi) \equiv \text{Model}^M(\pi)$$

All that remains is to bound collision probability.  
 Since  $r$  has empirical entropy  $\geq \lambda$  by def, this is at most  $T \cdot 2^{-\lambda} = \text{negl}(\lambda)$ .

# Completeness

Want:  $\Pr_{x \leftarrow \text{Wat}^0(\pi)} \left[ \text{Detect}^0(x) = \text{false} \text{ AND } H_e(M, \pi, x) \geq 8\lambda\sqrt{L} \right] \leq 2^{-\lambda}$ .

will show a stronger statement: for ANY  $x \in T^*$  and  $\pi$  st  $H_e(M, \pi, x) \geq 8\lambda\sqrt{L}$

& ~~each bit of  $x$  has empir.~~  
 $\Pr_{\pi} \left[ \text{Detect}^0(x) = \text{false} \mid \text{Model}^0(\pi) = x \right] \geq 2^{-\lambda}$ .

r.v.  $S_j = \ln\left(\frac{1}{V_j}\right)$  Conditioned on  $\text{Mark}^0(\Pi) = X$   
or eq. on  $X_j$ .

If  $X_j = 1$

$S_j = \ln\left(\frac{1}{U_j}\right)$  conditioned on  $U_j \leq p_j$   
 $\swarrow$  random in  $[0,1]$   $\searrow$  random in  $[0, p_j]$

So  $S_j$  is exponential conditioned  
on  $S_j \geq \ln \frac{1}{p_j}$

**Fact** exponential r.v.s are memoryless

So the conditional dist of  $S_j$  cond. on  $X_j = 1$  is

$\ln\left(\frac{1}{p_j}\right) + E_j$  when  $E_j$  is exponential  
with rate 1.

Similarly for  $X_j = 0$

.....  
the conditional dist of  $S_j$  cond. on  $X_j = 0$  is

$\ln\left(\frac{1}{1-p_j}\right) + E_j$  when  $E_j$  is exponential  
with rate 1.

$$\sum_{j=l+1}^L s_j = \ln(2) \cdot \underbrace{H_e^{[l+1:L]}(M, \pi, x)} + \sum_{j=l+1}^L \epsilon_j \quad (18)$$

$$\geq H_e(M, \pi, x) - O(\lambda)$$

$$\geq \cancel{\Omega(\lambda\sqrt{L})} \quad 2\lambda\sqrt{L-l}$$

$$\Pr \left[ \sum_{j=l+1}^L s_j \leq (L-l) + \lambda\sqrt{L-l} \right]$$

$$= (L-l) - \lambda\sqrt{L-l} + 2\lambda\sqrt{L-l}$$

$$\leq \Pr \left[ \sum_{j=l+1}^L \epsilon_j \leq (L-l) - \lambda\sqrt{L-l} \right]$$

$$\leq e^{-\lambda^2/2}$$

This comes  
from exponential  
r.v. tail bounds

$$\geq (L-l) - \lambda\sqrt{L-l}$$

$$\Rightarrow \text{~~amp~~ detection w.p.} \geq 1 - e^{-\lambda^2/2} \quad \blacksquare$$

~~~~~