

# Publicly Detectable Watermarking for Language Models

<https://arxiv.org/abs/2310.18491>

本质：把message - signature pair植入大模型生成的文本中

- 先生成一对secret, public key  $sk, pk$
- 大模型生成的前 $l$ 个词不受任何干预，这 $l$ 个词算出来的signature  $\sigma = \text{Sign}_{sk}(H(t[0:l]))$
- 设这个 $\sigma$ 有 $\lambda_\sigma$ 个bit，我们要依次把这 $\lambda_\sigma$ 个bit融入到后续的 $\lambda_\sigma$ 个chunk当中（每个chunk的size是 $l$ ），我们称这植入 $\sigma$ 的一整个过程叫一轮植入。
- 植入过程：

假设现在要植入 $\sigma$ 的第 $i$ 个bit，让大模型再sample出一段长度为 $l$ 的文本 $x$ ，计算 $H(m||x||\sigma_{prev})$ ，其中 $m$ 为这一轮植入中所有新生成的已经被接受的文本，长度应该为 $l \times (i - 1)$ ， $\sigma_{prev}$ 为已经被植入的signature片段，长度为 $i - 1$ ；若 $H(m||x||\sigma_{prev}) =$  余下的signature片段，则我们让 $x$ 被接受，否则重新生成 $x$ ；

- 反着解密，这里 $m$ 是暴露的，因为如果只考虑一个message-signature pair的植入的话， $m$ 就是除了前 $l$ 个token的文本。由 $m$ 和 $H$ 反着算出 $\sigma$ ，再验证 $m, \sigma$ 是否是一对。