

A Watermark for Large Language Models

<https://arxiv.org/pdf/2301.10226.pdf>

经典文章, Inference time Watermark奠基, Red/Green List范式; 把词表分成红/绿, 给绿词加bias

Soft Red list rule: Watermark generation & detection

- Analyze the generating process of the nth token



- Watermark Detection: Analyze the distribution of the number of green words

No watermark

Every position: γ green $(1-\gamma)$ red

Binomial distribution \rightarrow normal distribution

$$N(\gamma T, \gamma(1-\gamma)T)$$

$$\text{Evaluation Standard: } z = (|S|_G - \gamma T) / \sqrt{\gamma(1-\gamma)T}$$

For example: $z > 4 \rightarrow$ machine-produced!