

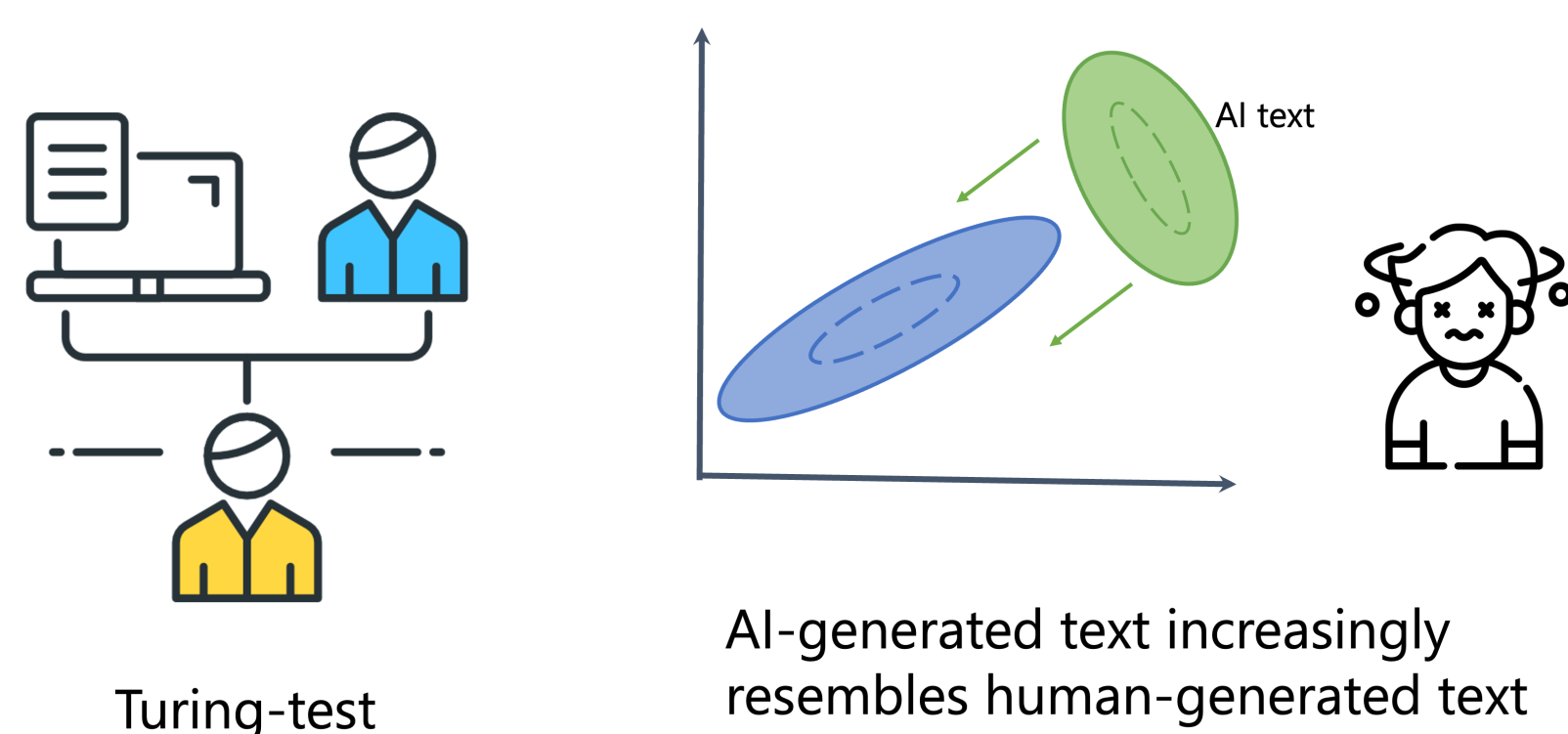


## Motivation

- **Potential harms of LLM**
  - Generate fake news
  - Contaminate web content
  - Assist in academic dishonesty
- **If most text in daily life is AI generated?**

Detect AI-generated text

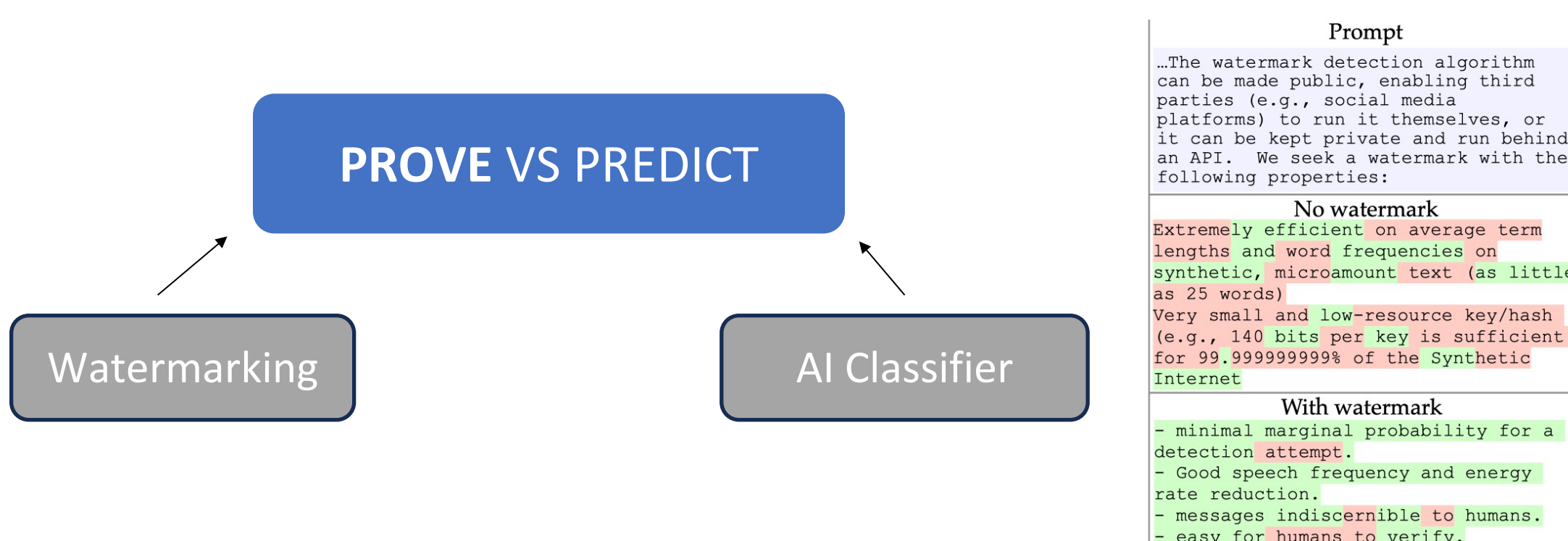
## Distinguish AI-generated text from human



## ZeroGPT/GPTZero/DetectGPT

- Not robust to distribution changes
- Prone to biases
- Vulnerable to adversarial attacks

## Watermarking digital text



## GPTWatermark

### Watermarking text generation

1. Randomly generate a watermark key  $k$ . Use watermark key to partition the vocabulary into a **Green List** of size  $\gamma|V|$
2. For  $t = 1, 2, \dots$ 
  1. Apply the language model to prior tokens to obtain a logit vector  $\ell_t$
  2. **Add  $\delta$  to each green list logit.** Apply the Softmax operator

$$\hat{\mathbf{p}}_t[v] = \begin{cases} \frac{\exp(\ell_t[v] + \delta)}{\sum_{i \in Red} \exp(\ell_t[i]) + \sum_{i \in Green} \exp(\ell_t[i] + \delta)}, & v \in Green \\ \frac{\exp(\ell_t[v])}{\sum_{i \in Red} \exp(\ell_t[i]) + \sum_{i \in Green} \exp(\ell_t[i] + \delta)}, & v \in Red. \end{cases}$$

3. Decode the next token using the watermarked distribution  $\hat{\mathbf{p}}_t$

### Watermarking text detection

1. Use the watermark detection key  $k$  to find the Green List
2. Tokenize the suspect text and calculate the number of green list tokens  $|y|_G = \sum_{t=1}^n \mathbf{1}(y_t \in G)$
3. Assume the null hypothesis is  $H_0$ : *The text sequence is generated with no knowledge of the green list rule.* Compute the z-statistic for this test:

$$z = (|y|_G - \gamma n) / \sqrt{n\gamma(1-\gamma)}$$

4. If  $z > \tau$ , the suspect text is watermarked

## Theoretical framework

### $\omega$ -Quality of watermarked output

$$D(\hat{\mathbf{p}}_t || \mathbf{p}_t) \leq \omega$$

### $\alpha$ -Type I error ("No false positives"):

$$\mathbb{P} \left[ \text{Detect}(k, \mathbf{y}) = 1 ; \begin{matrix} (\hat{\mathcal{M}}, k) \sim \text{Watermark}(\mathcal{M}) \\ \mathbf{y} \sim \mathcal{A}(\mathbf{x}, \text{aux}) \end{matrix} \right] \leq \alpha(\mathbf{x}, \mathcal{M}).$$

### $\beta$ -Type II error ("No false negatives"):

$$\mathbb{P} \left[ \text{Detect}(k, \mathbf{y}) = 0 ; \begin{matrix} (\hat{\mathcal{M}}, k) \sim \text{Watermark}(\mathcal{M}) \\ \mathbf{y} \sim \mathcal{M}(\mathbf{x}) \end{matrix} \right] \leq \beta(\mathbf{x}, \mathcal{M}).$$

### Security property

The adversary needs to make enough edits to evade detection.

## Robustness property

### GPTWatermark robustness to editing

Twice the robustness!

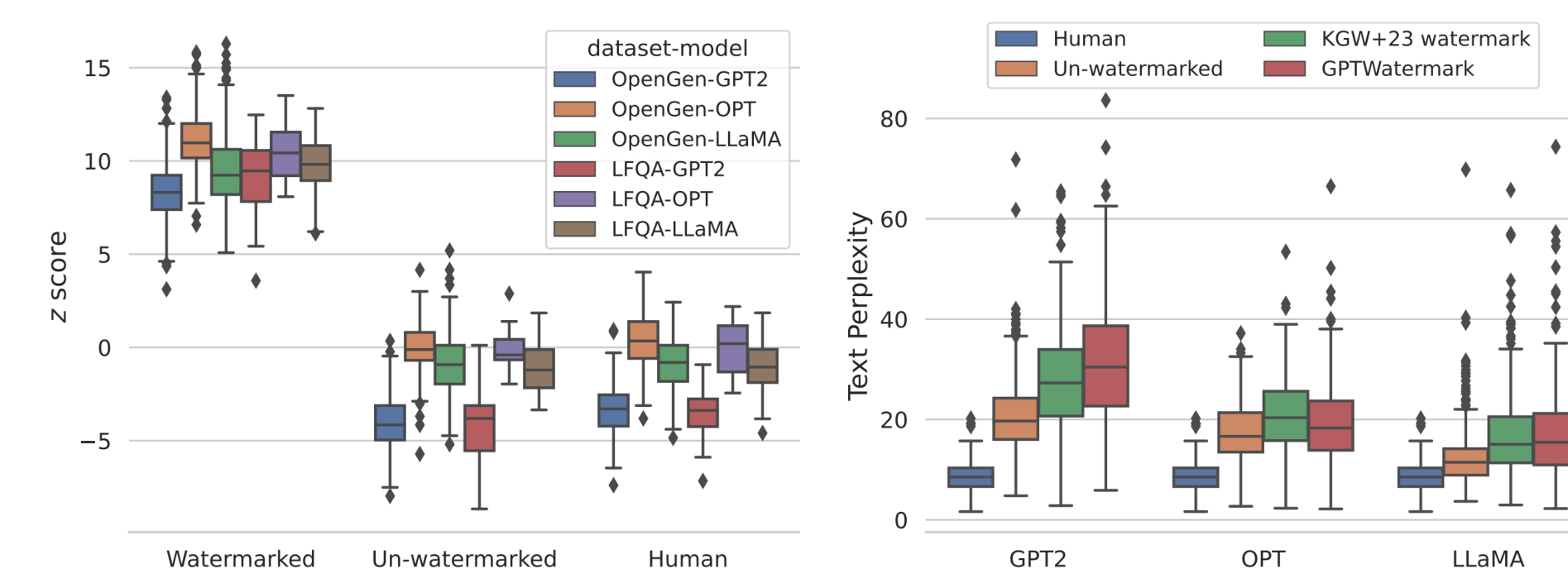
$$z_u \geq z_y - \max \left\{ \frac{(1 + \gamma/2)\eta}{\sqrt{n}}, \frac{(1 - \gamma/2)\eta}{\sqrt{n - \eta}} \right\}.$$

### KGW+23 watermark robustness to editing

$$z_u \geq z_y - \max \left\{ \frac{(2 + \gamma/2)\eta}{\sqrt{n}}, \frac{(2 - \gamma/2)\eta}{\sqrt{n - \eta}} \right\}.$$

## Experiment results

### z-score and text perplexity



### Robustness against paraphrasing attack

Setting	Method	OpenGen				LFQA			
		1% FPR		10% FPR		1% FPR		10% FPR	
		TPR	F1	TPR	F1	TPR	F1	TPR	F1
No attack	KGW+23	1.000	0.995	1.000	0.952	1.000	0.995	1.000	0.952
	GPTWatermark	1.000	0.995	1.000	0.952	1.000	0.995	1.000	0.952
ChatGPT	KGW+23	0.565	0.704	0.853	0.747	0.327	0.453	0.673	0.490
	GPTWatermark	0.866	0.910	0.961	0.818	0.442	0.568	0.865	0.584
DIPPER-1	KGW+23	0.386	0.546	0.738	0.720	0.372	0.534	0.740	0.767
	GPTWatermark	0.729	0.830	0.922	0.837	0.639	0.770	0.909	0.865
DIPPER-2	KGW+23	0.490	0.646	0.810	0.769	0.432	0.595	0.845	0.839
	GPTWatermark	0.777	0.862	0.941	0.852	0.693	0.810	0.948	0.894
BART	KGW+23	0.342	0.505	0.667	0.759	0.457	0.617	0.783	0.836
	GPTWatermark	0.590	0.730	0.861	0.857	0.656	0.784	0.885	0.897

### Distinguishing human-written text

