

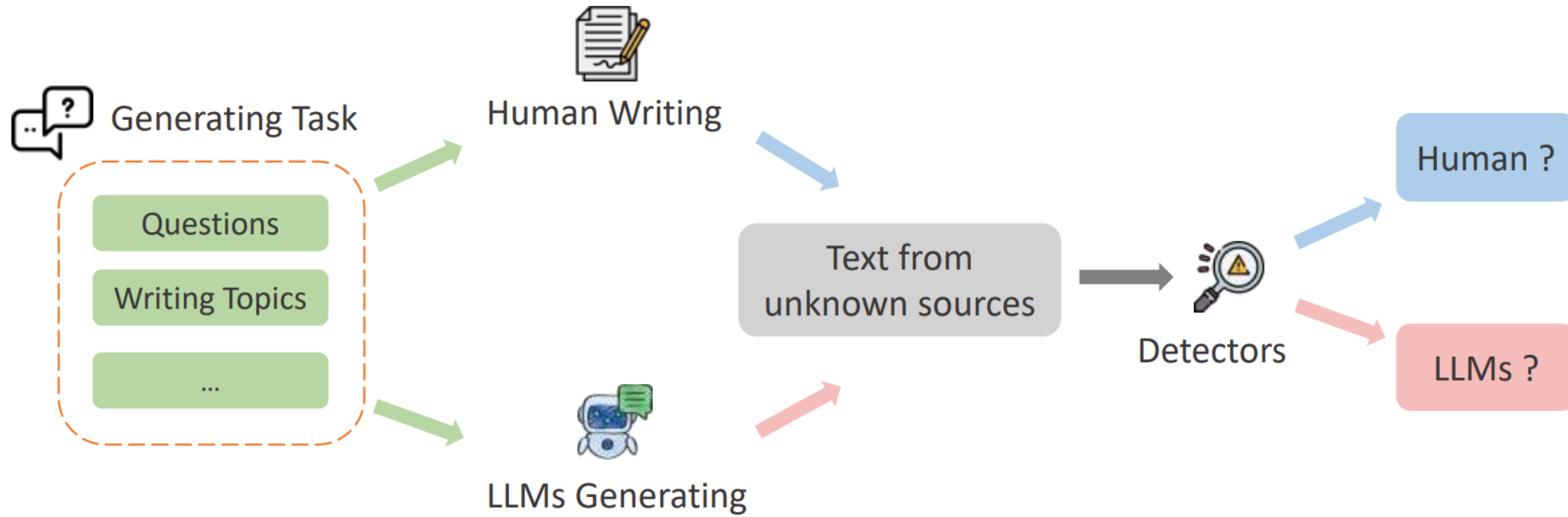
Ten Words Only Still Help: Improving Black-Box AI-Generated Text Detection via Proxy-Guided Efficient Re-Sampling

Yuhui Shi^{1,2}, Qiang Sheng¹, Juan Cao^{1,2}, Hao Mi^{1,2}, Beizhe Hu^{1,2}, Danding Wang¹

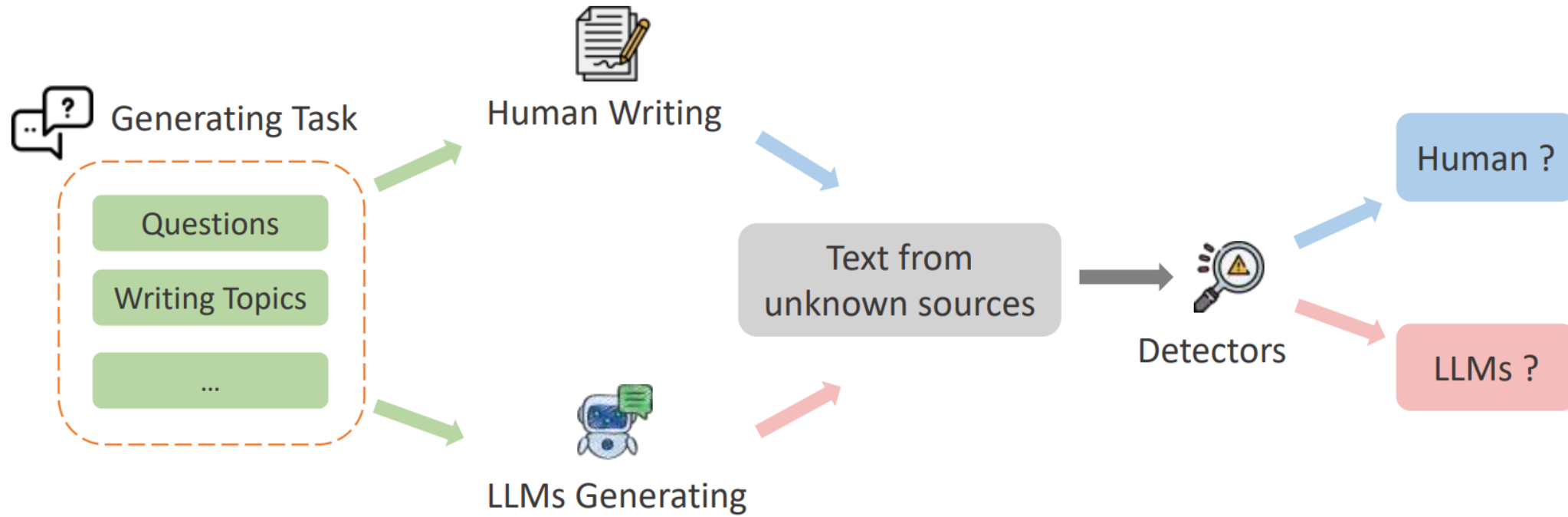
¹Institute of Computing Technology, Chinese Academy of Sciences

²University of Chinese Academy of Sciences

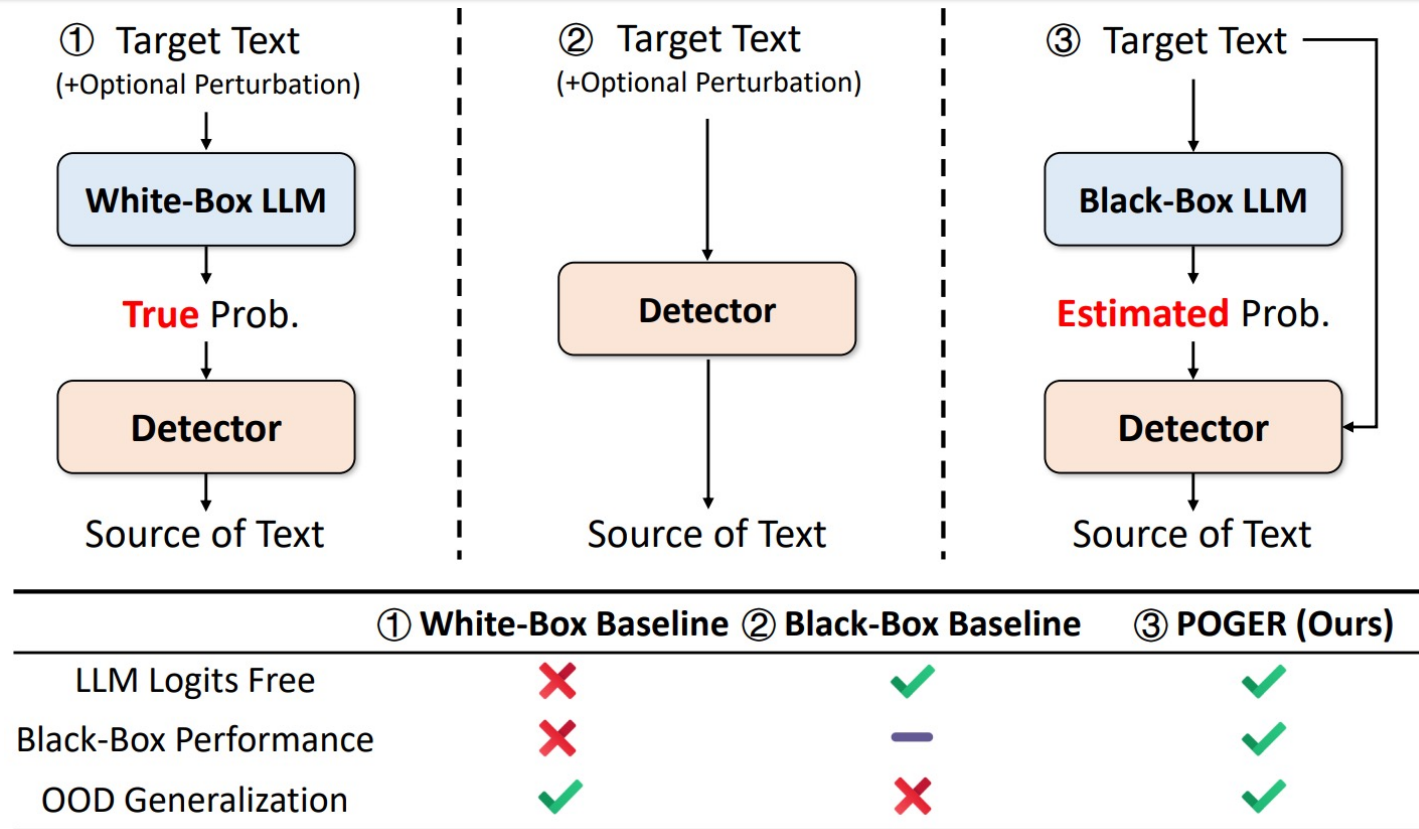
2024.08.09



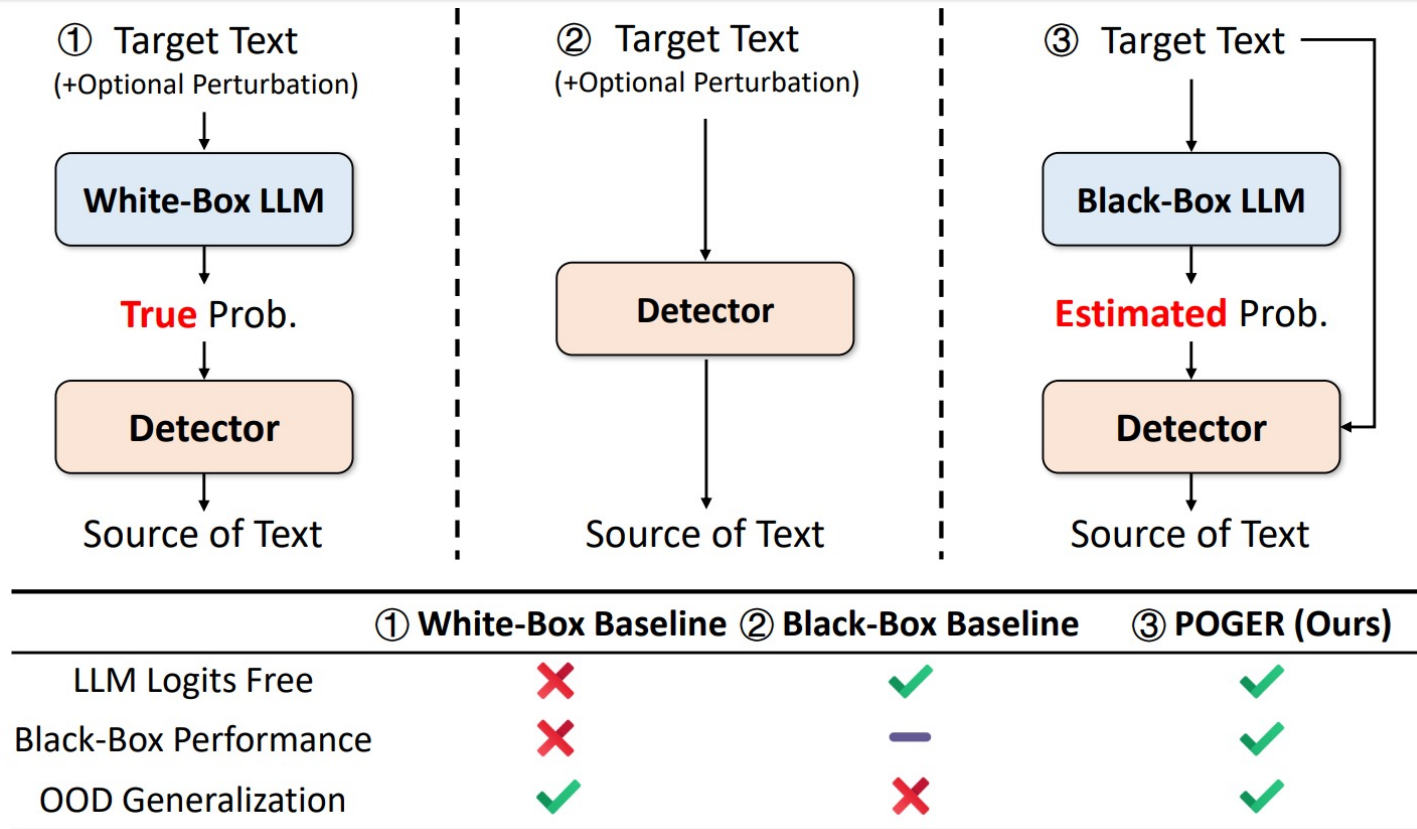
- **Background:** The misuse of large language models (LLMs) has led to issues such as misinformation and academic dishonesty, which makes **AI-generated text (AIGT) detection** critical.



- **Task Formulation:** AIGT detection aims to obtain a classifier $f: \mathbf{x} \rightarrow y$, where y is the source of the given text \mathbf{x} .
 - **Binary AIGT Detection:** $y \in \{\text{human, AI}\}$
 - **Multiclass AIGT Detection:** $y \in \{\text{human, } \theta_1, \theta_2, \dots, \theta_M\}$, where θ_i is a LLM



➤ **Challenge:** White-box methods have better performance and generalizability, but they require access to LLMs' internal states and are **not applicable to black-box settings**.



- **Solution:** Estimate word generation probabilities as pseudo white-box features via **multiple re-sampling** to help improve AGT detection under the black-box setting.

A naive solution:

- For each word in given text x , we instruct the black-box LLM for N times using the following prompt:

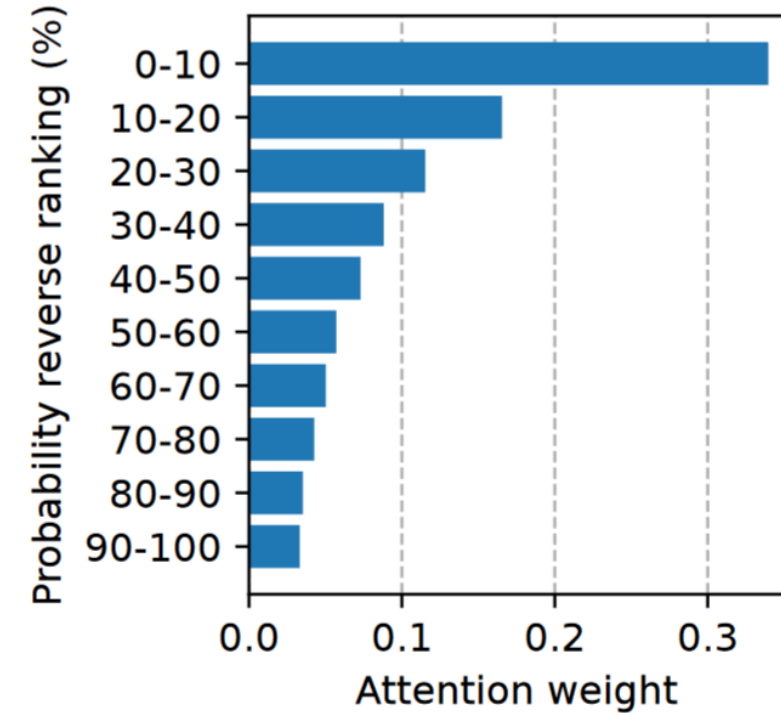
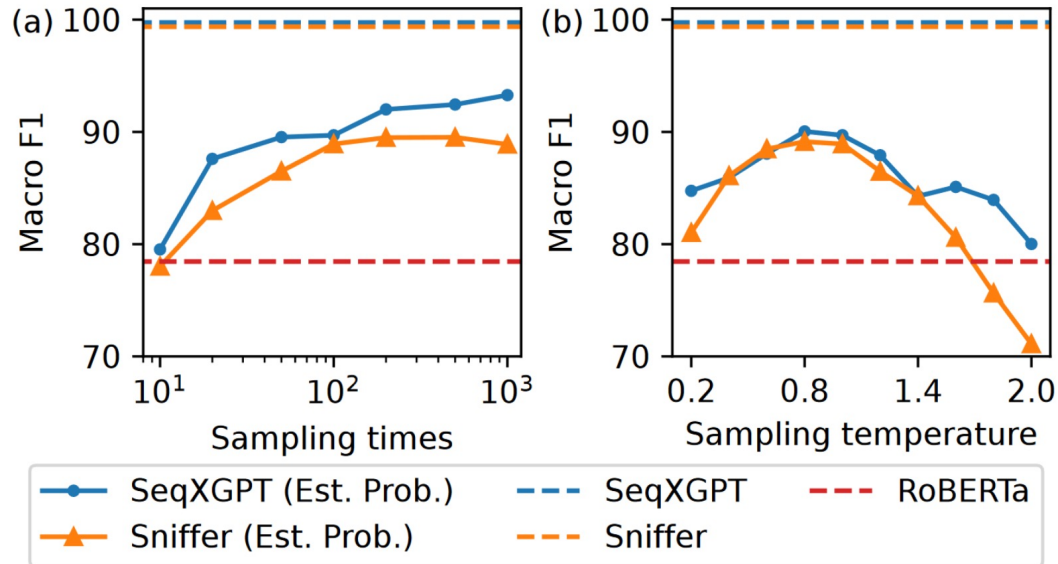
Please continue writing the following text, starting from the next word: $\{x_{<i}\}$.

- The estimated probability of x_i given $\{x_{<i}\}$ is computed as the frequency of x_i in the output word set $\{o_j\}_{j=1}^N$:

$$\hat{p}(x_i|x_{<i}) = \frac{1}{N} \sum_{j=1}^N \mathbb{I}(o_j = x_i).$$

- Use estimated probability as **an alternative input** of white-box methods.

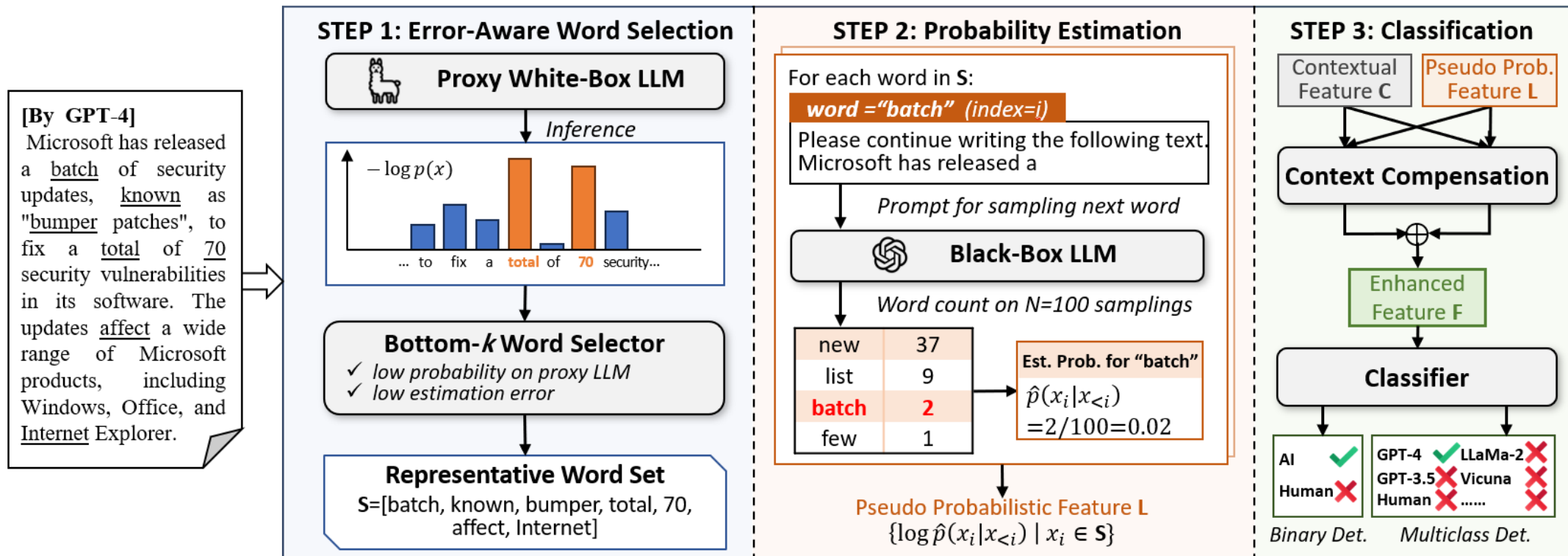
Preliminary Study



Finding 1: It is feasible to perform black-box AIGT detection by estimated probs.

Finding 2: **Low-probability words** gain higher attention from the detector.

We propose POGER, a **proxy-guided efficient re-sampling** method.

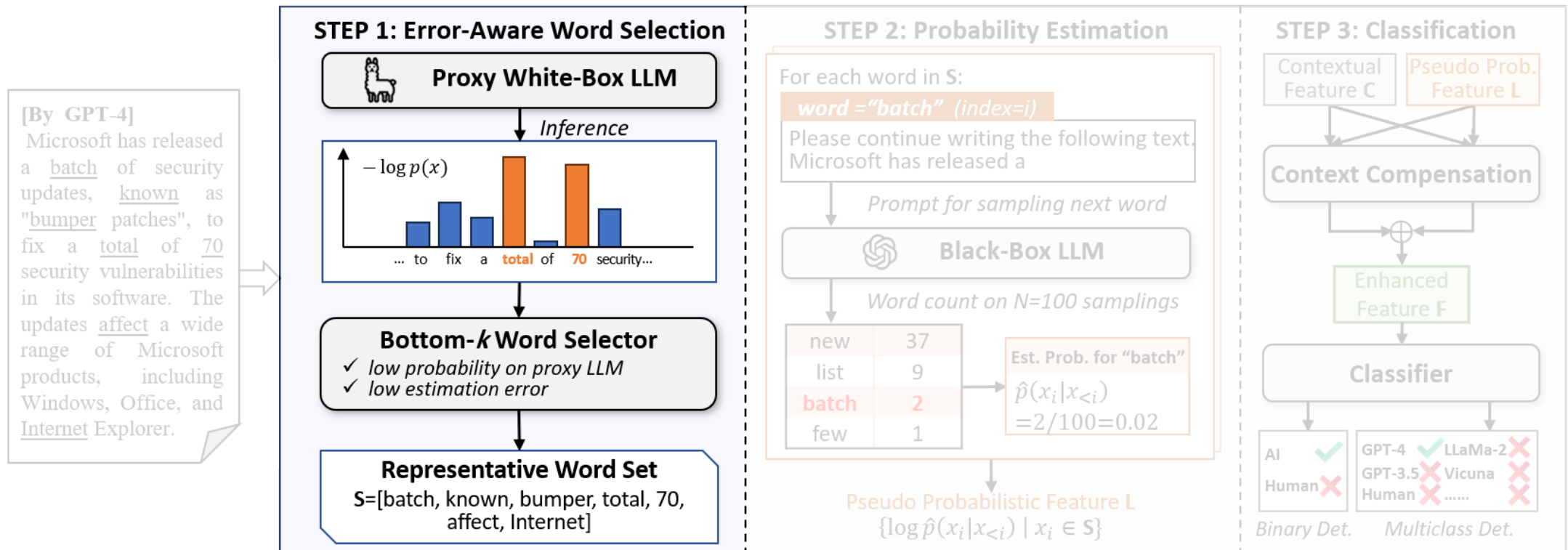


Step 1: Error-Aware Word Selection



- Use an easy-to-use LM (e.g., GPT-2) as the **proxy** to infer on the given text x and obtain token probabilities $p^\theta = (p_1^\theta, p_2^\theta, \dots, p_n^\theta)$
- Adopt an **error-aware** bottom- k word selector to get the representative word set S :

$$p^{\theta'} = \left\{ p_i \mid p_i \geq \frac{1}{1 + N\Delta^2} \right\} \quad \text{IDX} = \{i \mid p_i^\theta \in \text{MINK}(p^{\theta'})\}, S = \{x_i \mid i \in \text{IDX}\}$$

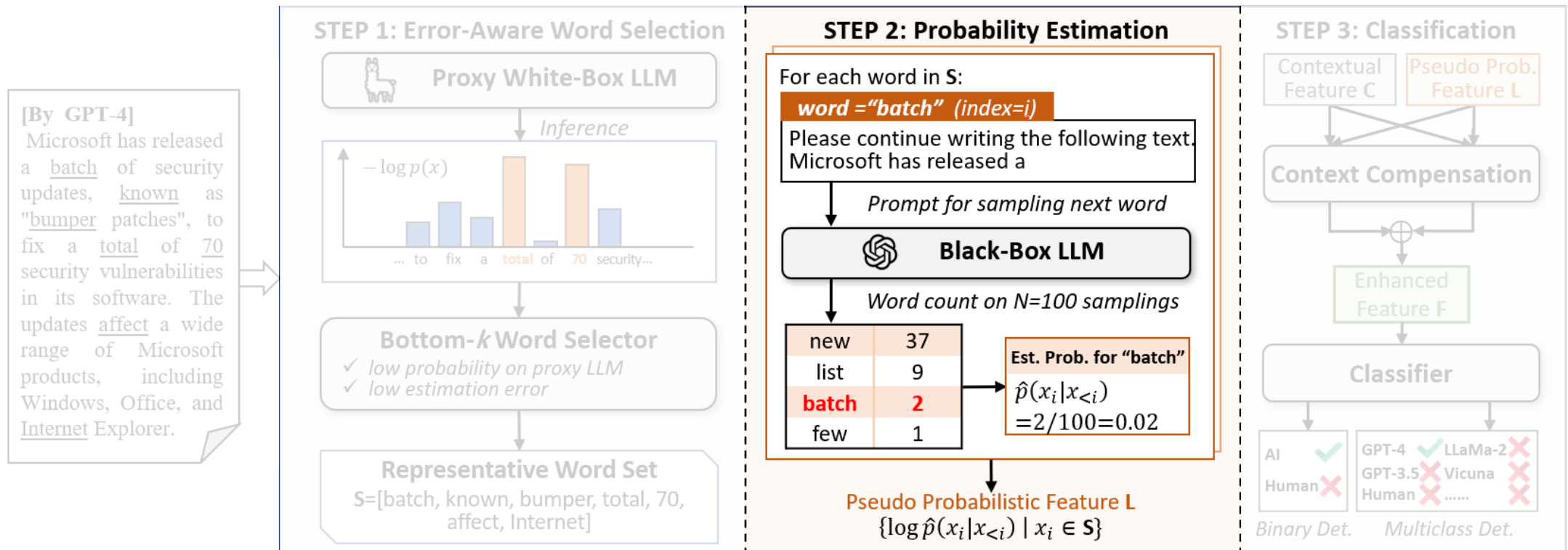


Step 2: Probability Estimation



- **Sample** and **calculate** probability for the selected k words in S on the given M candidate black-box LLMs by N times
- Get the **pseudo log probabilistic feature** matrix:

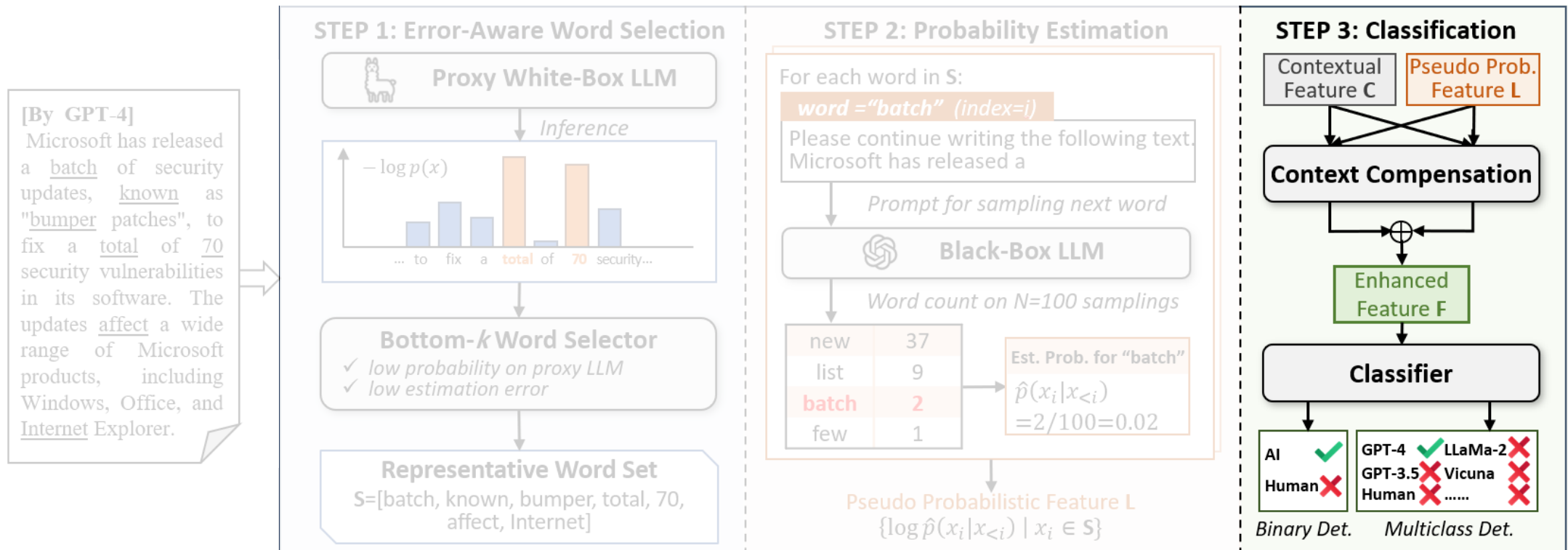
$$\mathbf{L} = [l_i]_{i=1}^k \in \mathbb{R}^{k \times M} \quad l_i = \left[\hat{p}_{\theta_j} (x_{\text{IDX}[i]} | x_{\text{IDX}[i]-b:\text{IDX}[i]-1}) \right]_{j=1}^M$$



Step 3: Context-Compensated Classification



- As **context compensation**, introduce the contextual semantic representation $\mathbf{C} \in \mathbb{R}^{k \times d}$
- $\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$, $\mathbf{F} = \text{Attention}(\mathbf{L}, \mathbf{C}, \mathbf{C}) \oplus \text{Attention}(\mathbf{C}, \mathbf{L}, \mathbf{L})$
- $\hat{y} = \text{softmax}(\text{MLP}(\mathbf{F}))$



Method	Human	GPT-2	GPT-J	LLaMA-2	Vicuna	Alpaca	GPT-3.5	GPT-4	MacF1
Partial White-Box Setting									
DNA-GPT White Sniffer	N/A	62.70	40.79	45.36	30.49	70.18	N/A	N/A	49.91*
SeqXGPT	96.60	100.00	100.00	<u>98.49</u>	95.85	99.23	75.34	72.65	92.27
POGER-Mixture	98.07	100.00	<u>99.62</u>	98.88	99.62	<u>98.87</u>	85.93	84.17	95.64
<i>w/o Context Compensation</i>	<u>97.32</u>	98.88	99.23	98.11	<u>97.71</u>	98.86	97.36	97.38	98.11
	96.97	<u>99.62</u>	99.23	96.68	94.94	98.48	<u>95.42</u>	<u>95.13</u>	<u>97.06</u>
Black-Box Setting									
RoBERTa	88.24	78.03	86.55	55.47	58.70	59.91	70.63	84.13	72.71
T5-Sentinel	87.29	85.42	<u>88.71</u>	67.78	62.11	69.73	75.79	79.83	77.08
DNA-GPT Black Sniffer	N/A	38.58	21.56	48.80	33.85	47.15	53.99	39.82	40.53*
SeqXGPT	87.41	<u>89.82</u>	87.26	29.52	47.62	35.84	34.21	52.63	58.04
POGER	<u>91.67</u>	89.66	86.77	23.64	46.31	45.64	42.10	62.40	61.02
<i>w/o Context Compensation</i>	92.49	93.75	89.96	90.49	89.30	93.82	90.98	92.59	91.67
	84.21	88.30	80.63	<u>81.88</u>	<u>88.65</u>	<u>91.95</u>	<u>89.49</u>	<u>87.35</u>	<u>86.56</u>

- POGER **outperforms all baseline methods** in both settings of multiclass AIGT detection.
- POGER has **better OOD generalization capabilities**, benefiting from the pseudo probabilistic.

Method	In-Dist.	Out-of-Distribution			
		QA→Writing		Writing→QA	
RoBERTa	72.71	54.23	(-25.42%)	46.73	(-35.73%)
T5-Sentinel	77.08	47.23	(-38.73%)	53.19	(-30.99%)
Sniffer	58.04	57.50	(-0.93%)	53.16	(-8.41%)
SeqXGPT	61.02	59.07	(-3.20%)	54.94	(-9.96%)
POGER	91.67	89.00	(-2.91%)	84.19	(-8.16%)

- **Motivation:** Achieve “**white-boxing**” the black-box LLM by estimating word generation probability through multiple re-sampling, so that the high performance white-box detection methods can also be used under black-box setting.
- **Method:** By selecting **low-probability words** as representative words, the number of re-samples can be significantly reduced, thus improving efficiency and reducing costs.
- **Result:** Experiments on texts from humans and 7 LLMs demonstrated the superiority of POGER.

Data & Code



<https://github.com/ICTMCG/POGER>

Paper List



[https://github.com/ICTMCG/
Awesome-Machine-Generated-Text](https://github.com/ICTMCG/Awesome-Machine-Generated-Text)

Our Latest Work



(WeChat, in chinese)

THANKS

shiyuhui22s@ict.ac.cn