

# Efficient Detection of LLM-generated Texts with a Bayesian Surrogate Model(ACL 2024)

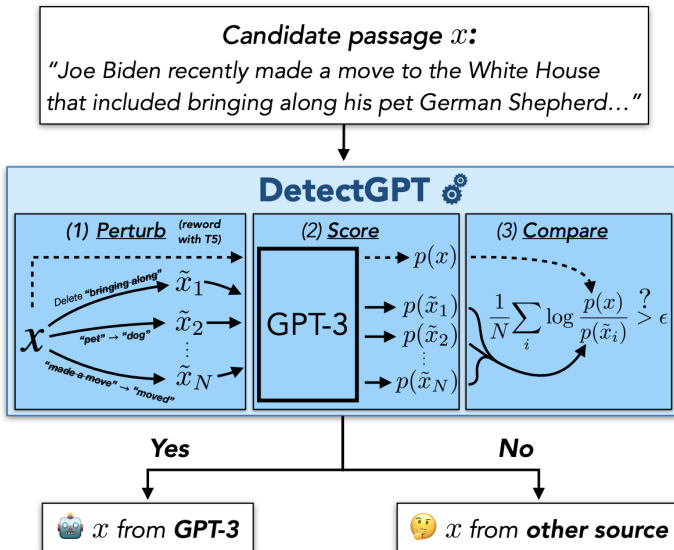
---

Kwan Ho Lee

August 4, 2025

Seoul National University

# DetectGPT



## Surrogate Model

- When an outcome of interest cannot be easily measured or computed, an approximate mathematical model of the outcome is used instead.
- In our case, we will fit a Gaussian Process model with the text as  $x$  and  $y = \log(p(x))$ . In other words, we aim to build a model that replaces steps (1) through (2) outlined above.

# Bayesian Surrogate Model

## Fitting Gaussian Process Regression

- Kernel function  $\mathcal{K}(x, x') = \alpha \cdot \text{BERTScore}(x, x') + \beta$
- Train Data
  - Make perturbation text  $\{x_i\}_{i=1}^N = \{x_1\} \cup X^*$
  - train with  $X_t = \{x_0(\text{machine-generated text}), x_1\}$
  - retrain with  $X_{t+1} = X_t \cup x \quad : x = \arg \max_{x \in X^*} \sigma_t^2(x)$
- Posterior Equations

$$p(f_{x^*} \mid X_t, y_t, X^*) = \mathcal{N}(\bar{f}^*, \Sigma^*)$$

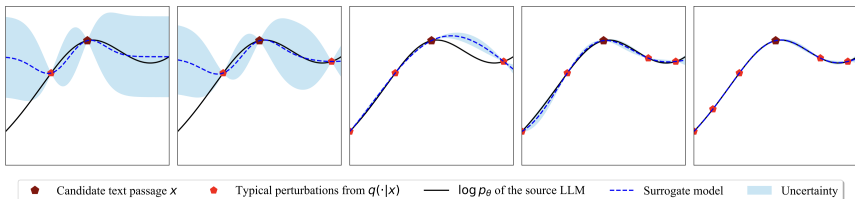
where

$$f(x) \sim \mathcal{GP}(0, k(x, x')) \quad (1)$$

$$\bar{f}^* := k_{x^*, x_t} [k_{x_t, x_t} + \sigma^2 I]^{-1} y_t \quad (2)$$

$$\Sigma^* := k_{x^*, x^*} - k_{x^*, x_t} [k_{x_t, x_t} + \sigma^2 I]^{-1} k_{x_t, x^*} \quad (3)$$

# Bayesian Surrogate Model



- Why choose a Gaussian Process as the surrogate model?
  - Performs well in low-data regimes (only a few queried points).
  - Resistance for overfitting
  - Capability to quantify uncertainty
  - Flexible kernels capture non-trivial local curvature.
- Query efficiency: This method detects a single machine-generated text with just 4 queries, whereas DetectGPT requires about 100 for the same task.

# Experiment

---

## Experimenting with different text-generation models

Dataset	Method	2	3	4	10	20	50	100	200
Xsum	DetectGPT (7B)	0.850	0.893	0.904	0.927	0.932	0.95	0.952	0.952
	<b>Our Method (7B)</b>	0.958	0.955	0.957	-	-	-	-	-
	DetectGPT (13B)	0.817	0.849	0.861	0.913	0.922	0.936	0.932	0.935
	<b>Our Method (13B)</b>	0.886	0.912	0.929	-	-	-	-	-
Squad	DetectGPT (7B)	0.798	0.820	0.857	0.871	0.878	0.889	0.886	0.884
	<b>Our Method (7B)</b>	0.947	0.936	0.932	-	-	-	-	-
	DetectGPT (13B)	0.671	0.686	0.712	0.740	0.731	0.743	0.758	0.757
	<b>Our Method (13B)</b>	0.799	0.785	0.787	-	-	-	-	-
Writing	DetectGPT (7B)	0.903	0.916	0.930	0.938	0.937	0.936	0.941	0.938
	<b>Our Method (7B)</b>	0.985	0.983	0.979	-	-	-	-	-
	DetectGPT (13B)	0.856	0.897	0.933	0.962	0.963	0.963	0.966	0.967
	<b>Our Method (13B)</b>	0.929	0.971	0.979	-	-	-	-	-

**Figure 1:** The AUROC for detecting samples generated by Vicuna-7B/13B varies depending on the number of queries made to the source model.

## White Box Assumption

Method	Xsum	Squad	Writing
DetectGPT (LLaMA2)	0.397	0.473	0.641
<b>Our Method (LLaMA2)</b>	0.631	0.660	0.742
DetectGPT (Vicuna)	0.595	0.606	0.733
<b>Our Method (Vicuna)</b>	0.780	0.714	0.850

**Figure 2:** The AUROC for detecting texts generated by ChatGPT with query budget 15



# Appendix

---

# Appendix - Number of Typical Samples

## The Visualization Typical Samples

Row Mean	BERT Score	Text
0.9227	0.92	Joe Biden recently made a move to the White House.
0.9164	0.930.90	The White House is now under the leadership of Joe Biden.
0.9257	0.920.950.91	Joe Biden has made his official residence at the White House.
0.9236	0.940.920.920.91	The White House is now the workplace of Joe Biden.
0.9282	0.940.910.930.920.93	Joe Biden has assumed the role of the White House occupant.
0.9272	0.940.920.940.910.930.92	Joe Biden's new address is the White House.
0.9313	0.940.920.930.920.940.930.94	Joe Biden has recently started his tenure at the White House.
0.9271	0.920.940.940.950.920.930.910.91	Joe Biden has established his administration in the White House.
0.9278	0.950.900.950.910.920.920.950.920.92	The White House is now Joe Biden's official residence.
0.9310	0.950.920.920.910.970.930.940.940.920.92	Joe Biden has recently made his way to the White House residence.
0.9313	0.940.910.950.920.920.940.950.920.920.960.92	Joe Biden has assumed the duties of the White House.
		Joe Biden has recently made the White House his new home.

**Figure 3:** The visualization of the candidate text passage and the first 11 typical perturbations of it identified by our method, ordered from top to bottom. The BertScore among them and the row mean (estimated without the diagonal elements) are reported. The log probabilities are given by GPT-2.

# References

- [1] Miao, Y., Gao, H., Zhang, H., Deng, Z. (2023). Efficient detection of LLM-generated texts with a Bayesian surrogate model. arXiv preprint arXiv:2305.16617.
- [2] Mitchell, E., Lee, Y., Khazatsky, A., Manning, C. D., Finn, C. (2023, July). Detectgpt: Zero-shot machine-generated text detection using probability curvature. In International conference on machine learning (pp. 24950-24962). PMLR.