



Duke

PRATT SCHOOL of  
ENGINEERING

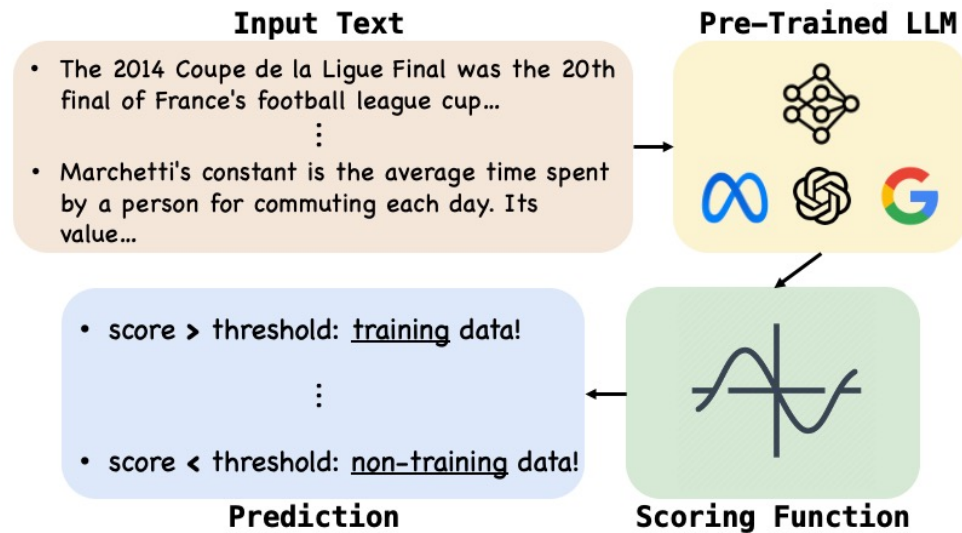
# Min-K%++: Improved Baseline for Detecting Pre-Training Data from Large Language Models

Jingyang Zhang<sup>1\*</sup>, Jingwei Sun<sup>1\*</sup>, Eric Yeats<sup>1</sup>, Yang Ouyang<sup>1</sup>,  
Martin Kuo<sup>1</sup>, Jianyi Zhang<sup>1</sup>, Hao Yang<sup>1,2</sup>, Hai Li<sup>1</sup>

<sup>1</sup>Duke University, <sup>2</sup>Johns Hopkins University



# Pre-training data detection for LLMs



Given a data instance  $x$  and the (gray-box) access to a pre-trained model  $\mathcal{M}$ , we want to give binary prediction of whether the sample  $x$  appears in the training set  $\mathcal{D}$ , with a scoring function  $s$ :

$$\text{prediction}(x, \mathcal{M}) = \begin{cases} 1 & (x \in \mathcal{D}), \quad s(x; \mathcal{M}) \geq \lambda \\ 0 & (x \notin \mathcal{D}), \quad s(x; \mathcal{M}) < \lambda \end{cases}$$

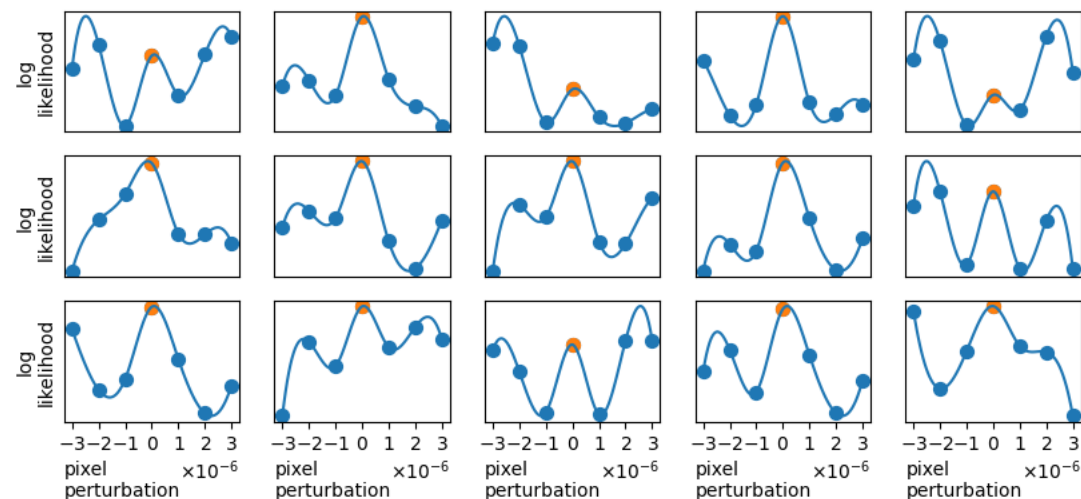
- This is an instance of Membership Inference Attack (MIA) that explicitly focuses on pre-trained LLMs, which poses unique challenges due to LLM pre-training setup.
- This problem has implications for critical issues such as copyright violation, test data contamination, and privacy auditing.

# Revisiting the training objective

- LLM does maximum likelihood training, which is equivalent to *implicit score matching* (ISM).
- The minimization of ISM indicates that **training samples tend to form local maxima along each input dimension** in the likelihood space.
- This insight can be empirically confirmed by examining a pre-trained diffusion model.

$$\frac{1}{N} \sum_{\mathbf{x}} \left[ \frac{1}{2} \|\psi(\mathbf{x})\|^2 + \underbrace{\sum_{i=1}^d \frac{\partial \psi_i(\mathbf{x})}{\partial x_i}}_{\text{the sum of the second-order partial derivatives}} \right],$$

$$\psi(\mathbf{x}) = \frac{\partial \log p(\mathbf{x})}{\partial \mathbf{x}}$$



## Translating the insight to a practical method

$$\text{Min-K}\%++_{\text{token seq.}}(x_{<t}, x_t) = \frac{\log p(x_t|x_{<t}) - \mu_{\cdot|x_{<t}}}{\sigma_{\cdot|x_{<t}}},$$

$$\text{Min-K}\%++(x) = \frac{1}{|\text{min-k}\%|} \sum_{(x_{<t}, x_t) \in \text{min-k}\%} \text{Min-K}\%++_{\text{token seq.}}(x_{<t}, x_t).$$

$$\mu_{x_{<t}} = \mathbb{E}_{z \sim p(\cdot|x_{<t})}[\log p(z|x_{<t})] \quad \sigma_{x_{<t}} = \sqrt{\mathbb{E}_{z \sim p(\cdot|x_{<t})}[(\log p(z|x_{<t}) - \mu_{x_{<t}})^2]}$$

- Analogous to detecting local maximum in continuous distribution, we propose to identify whether each token of the input sequence has larger probability than other candidates.
- Min-K%++ introduces no computation overhead and requires no extra model for detection.

# Improved performance on two benchmarks

Table 1: AUROC results on WikiMIA benchmark (Shi et al., 2024). *Ori.* and *Para.* denote the original and paraphrased settings, respectively. **Bolded** number shows the best result within each column across all methods. The proposed Min-K%++ leads to remarkable improvements over existing methods in most settings.

Len.	Method	Mamba-1.4B		Pythia-6.9B		LLaMA-13B		LLaMA-30B		LLaMA-65B		Average	
		<i>Ori.</i>	<i>Para.</i>	<i>Ori.</i>	<i>Para.</i>	<i>Ori.</i>	<i>Para.</i>	<i>Ori.</i>	<i>Para.</i>	<i>Ori.</i>	<i>Para.</i>	<i>Ori.</i>	<i>Para.</i>
32	Loss	61.0	61.4	63.8	64.1	67.5	68.0	69.4	70.2	70.7	71.8	66.5	67.1
	Ref	62.2	62.3	63.6	63.5	57.9	56.2	63.5	62.4	68.8	68.2	63.2	62.5
	Lowercase	60.9	60.6	62.2	61.7	64.0	63.2	64.1	61.2	66.5	64.8	63.5	62.3
	Zlib	61.9	62.3	64.3	64.2	67.8	68.3	69.8	70.4	71.1	72.0	67.0	67.4
	Neighbor	64.1	63.6	65.8	65.5	65.8	65.0	67.6	66.3	69.6	68.7	66.6	65.8
	Min-K%	63.2	62.9	66.3	65.2	68.0	68.4	70.1	70.7	71.3	72.2	67.8	67.9
	Min-K%++	<b>66.8</b>	<b>66.1</b>	<b>70.3</b>	<b>68.0</b>	<b>84.8</b>	<b>82.7</b>	<b>84.3</b>	<b>81.2</b>	<b>85.1</b>	<b>81.4</b>	<b>78.3</b>	<b>75.9</b>
64	Loss	58.2	56.4	60.7	59.3	63.6	63.1	66.2	65.5	67.9	67.7	63.3	62.4
	Ref	60.6	59.6	62.4	62.9	63.4	60.9	69.0	65.4	73.4	71.0	65.8	63.9
	Lowercase	57.0	57.0	58.2	57.7	62.0	61.0	62.1	59.8	64.5	61.9	60.8	59.5
	Zlib	60.4	59.1	62.6	61.6	65.3	65.3	67.5	67.4	69.1	69.3	65.0	64.5
	Neighbor	60.6	60.6	63.2	63.1	64.1	64.7	67.1	66.7	69.6	69.5	64.9	64.9
	Min-K%	62.2	58.0	65.0	61.1	66.0	64.0	68.5	65.7	69.8	67.9	66.3	63.3
	Min-K%++	<b>67.2</b>	<b>63.3</b>	<b>71.6</b>	<b>64.8</b>	<b>85.7</b>	<b>78.8</b>	<b>84.7</b>	<b>74.9</b>	<b>83.8</b>	<b>74.0</b>	<b>78.6</b>	<b>71.2</b>
128	Loss	63.3	62.7	65.1	64.7	67.8	67.2	70.3	69.2	70.7	70.2	67.4	66.8
	Ref	62.0	61.1	63.3	62.9	62.6	59.7	71.9	70.0	73.7	72.0	66.7	65.1
	Lowercase	58.5	57.7	60.5	60.0	60.6	56.4	59.1	55.4	63.3	60.1	60.4	57.9
	Zlib	65.6	65.3	67.6	<b>67.4</b>	69.7	69.6	71.8	71.5	72.1	<b>72.1</b>	69.4	69.2
	Neighbor	64.8	62.6	67.5	64.3	68.3	64.0	72.2	67.2	73.7	70.3	69.3	65.7
	Min-K%	66.8	64.5	69.5	67.0	71.5	68.7	73.9	70.2	73.6	70.8	71.0	68.2
	Min-K%++	<b>68.8</b>	<b>65.6</b>	<b>70.7</b>	66.8	<b>83.9</b>	<b>76.2</b>	<b>82.6</b>	<b>73.8</b>	<b>80.0</b>	70.7	<b>77.2</b>	<b>70.6</b>

Table 2: AUROC results on the challenging MIMIR benchmark (Duan et al., 2024) with a suite of Pythia models (Biderman et al., 2023). In each column, the best result across all methods is **bolded**, with the runner-up underlined. <sup>†</sup>Ref relies on an extra reference LLM. <sup>‡</sup>Neighbor induces significant extra computational cost than others (25× in this case), for which reason we don’t run on the 12B model. Despite not requiring a reference model like the Ref method does, our Min-K%++ consistently achieves superior or comparable performance.

Method	Wikipedia					Github					Pile CC					PubMed Central				
	160M	1.4B	2.8B	6.9B	12B	160M	1.4B	2.8B	6.9B	12B	160M	1.4B	2.8B	6.9B	12B	160M	1.4B	2.8B	6.9B	12B
Loss	50.2	51.3	51.8	52.8	53.5	65.7	69.8	71.3	73.0	74.0	49.6	50.0	50.1	50.7	51.1	49.9	49.8	49.9	50.6	51.3
<sup>†</sup> Ref	<b>51.2</b>	<b>55.2</b>	<b>58.1</b>	<b>61.8</b>	<b>63.9</b>	63.9	67.1	65.3	64.4	63.0	49.2	<b>52.2</b>	<b>53.7</b>	<b>54.9</b>	<b>56.7</b>	<b>51.3</b>	<b>53.1</b>	<b>53.7</b>	<b>55.9</b>	<b>58.2</b>
Zlib	<u>51.1</u>	52.0	52.4	53.5	54.3	<b>67.4</b>	<b>71.0</b>	<b>72.3</b>	<b>73.9</b>	<b>74.8</b>	49.6	50.1	50.3	50.8	51.1	49.9	50.0	50.1	50.6	51.2
<sup>‡</sup> Neighbor	50.7	51.7	52.2	53.2	/	65.3	69.4	70.5	72.1	/	49.6	50.0	50.1	50.8	/	47.9	49.1	49.7	50.1	/
Min-K%	50.2	51.3	51.8	53.6	54.4	<u>65.7</u>	69.9	71.4	<u>73.2</u>	<u>74.3</u>	<u>50.3</u>	51.0	50.8	51.5	51.7	50.6	50.3	50.5	51.2	52.3
Min-K%++	49.7	<u>53.7</u>	<u>55.1</u>	<u>58.0</u>	<u>61.1</u>	64.8	69.6	70.9	72.8	74.2	<b>50.6</b>	<u>51.0</u>	<u>51.0</u>	<u>53.0</u>	<u>53.5</u>	<u>50.6</u>	<u>51.4</u>	<u>52.4</u>	<u>54.2</u>	<u>55.4</u>

Method	ArXiv					DM Mathematics					HackerNews					Average				
	160M	1.4B	2.8B	6.9B	12B	160M	1.4B	2.8B	6.9B	12B	160M	1.4B	2.8B	6.9B	12B	160M	1.4B	2.8B	6.9B	12B
Loss	<u>51.0</u>	51.5	51.9	52.9	53.4	48.8	48.5	48.4	48.5	48.5	49.4	50.5	51.3	52.1	52.8	52.1	53.1	53.5	54.4	54.9
<sup>†</sup> Ref	49.4	<u>51.5</u>	<u>53.1</u>	<b>55.8</b>	<u>57.5</u>	<b>51.1</b>	<b>51.1</b>	<u>50.5</u>	<u>51.1</u>	<u>50.9</u>	49.1	<b>52.2</b>	<b>55.1</b>	<b>57.9</b>	<b>60.6</b>	52.2	<b>54.6</b>	<b>55.6</b>	<b>57.4</b>	<u>58.7</u>
Zlib	50.1	50.9	51.3	52.2	52.7	48.1	48.2	48.0	48.1	48.1	49.7	50.3	50.8	51.2	51.7	52.3	53.2	53.6	54.3	54.8
<sup>‡</sup> Neighbor	50.7	51.4	51.8	52.2	/	49.0	47.0	46.8	46.6	/	<u>50.9</u>	<u>51.7</u>	51.5	51.9	/	52.0	52.9	53.2	53.8	/
Min-K%	<b>51.0</b>	<b>51.7</b>	52.5	53.6	54.6	49.4	49.7	49.5	49.6	49.7	<b>50.9</b>	51.3	52.6	53.6	54.6	<b>52.6</b>	53.6	54.2	55.2	55.9
Min-K%++	50.1	51.1	<b>53.7</b>	<u>55.2</u>	<b>58.0</b>	<u>50.5</u>	<u>50.9</u>	<b>51.7</b>	<b>51.6</b>	<b>51.9</b>	50.7	51.3	<u>52.6</u>	<u>54.5</u>	<u>56.5</u>	<u>52.4</u>	<u>54.1</u>	<u>55.3</u>	<u>57.0</u>	<b>58.7</b>

- Min-K%++ consistently exhibits improved detection rates over existing reference-free methods.
- Compared with reference-based methods, Min-K%++ is comparable and superior.

# Conclusions

- We provide a theoretical insight that translates the training data detection problem into the identification of local maxima.
- We propose Min-K%++, a novel and well-performing method for LLM pre-training data detection.
- We demonstrate improved detection rates on two established benchmarks, setting a new baseline for the field.

Code and more information:

<https://github.com/zjysteven/mink-plus-plus>

