

# How Stable Is The Next Token? A Geometric View Of LLM Prediction Stability

Deyuan Liu<sup>1</sup>, Zecheng Wang<sup>1,2</sup>, Zhanyue Qin<sup>1</sup>, Zhiying Tu<sup>1</sup>, Dianhui Chu<sup>1</sup>, Dianbo Sui<sup>1\*</sup>

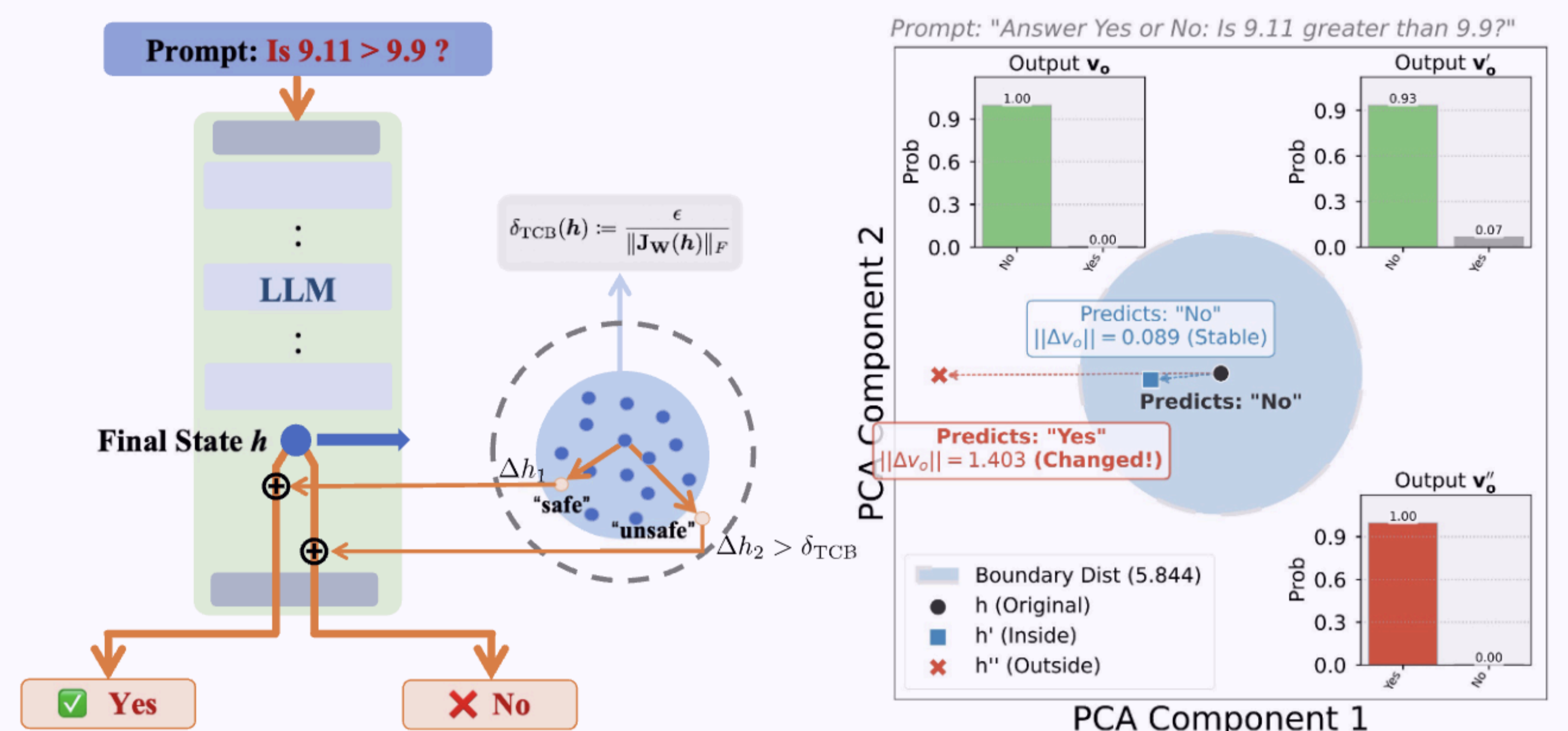
<sup>1</sup>Harbin Institute of Technology, <sup>2</sup>Wechat AI

deyuanliu@stu.hit.edu.cn, suidianbo@hit.edu.cn



## Introduction

Large Language Models (LLMs) exhibit impressive capabilities yet suffer from sensitivity to slight input context variations, hampering reliability. Conventional metrics like accuracy and perplexity fail to assess local prediction robustness, as normalized output probabilities can obscure the underlying resilience of an LLM's internal state to perturbations. We introduce the **Token Constraint Bound** ( $\delta_{TCB}$ ), a novel metric that quantifies the maximum internal state perturbation an LLM can withstand before its dominant next-token prediction significantly changes. Intrinsically linked to output embedding space geometry,  $\delta_{TCB}$  provides insights into the stability of the model's internal predictive commitment. Our experiments show  $\delta_{TCB}$  correlates with effective prompt engineering and uncovers critical prediction instabilities missed by perplexity during in-context learning and text generation.  $\delta_{TCB}$  offers a principled, complementary approach to analyze and potentially improve the contextual stability of LLM predictions.



## Definition

**Definition 1 (Token Constraint Bound  $\delta_{TCB}$ ).** Given the output weight matrix  $\mathbf{W}$ , hidden state  $\mathbf{h}$ , resulting output distribution  $\mathbf{o} = \text{softmax}(\mathbf{W}\mathbf{h})$ , and a tolerance  $\epsilon > 0$  for the maximum  $L_2$  change allowed in  $\mathbf{o}$ , the Token Constraint Bound  $\delta_{TCB}$  at state  $\mathbf{h}$  is defined as:

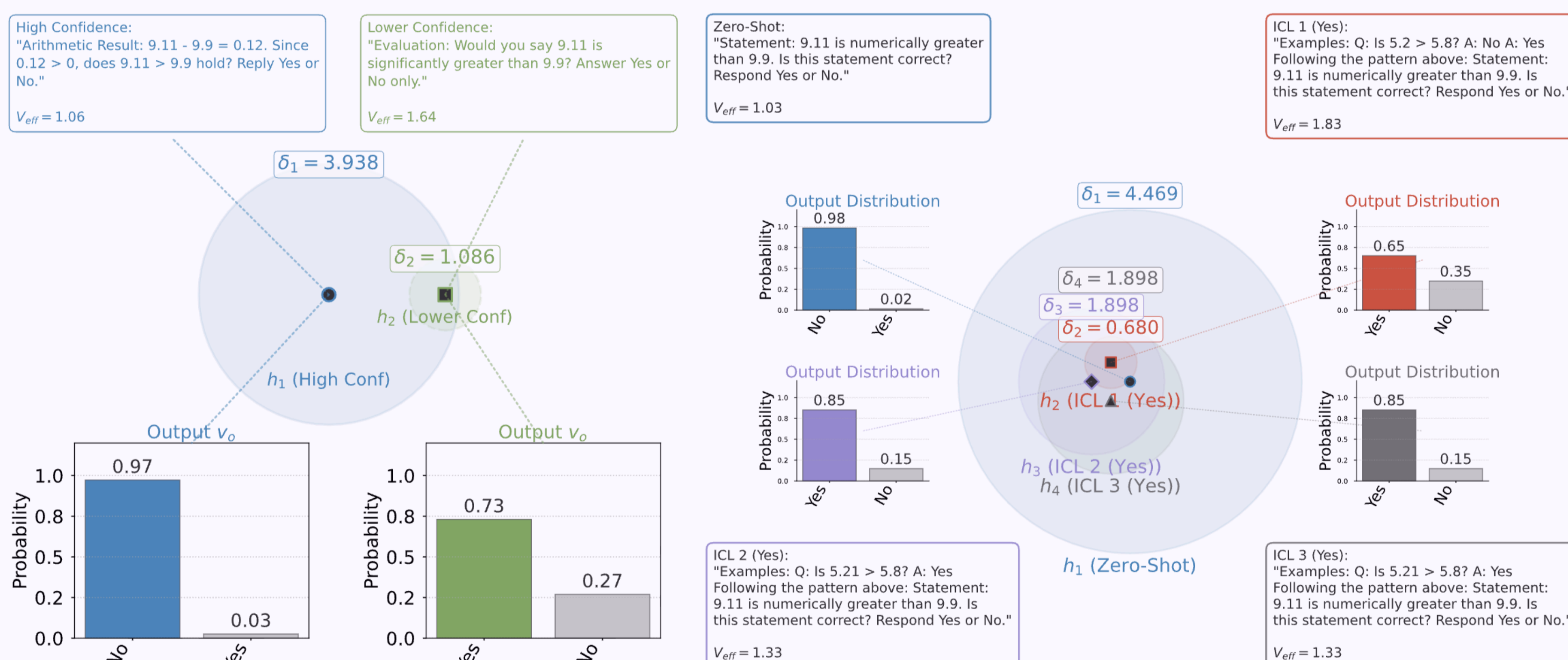
$$\delta_{TCB}(\mathbf{h}) := \frac{\epsilon}{\|\mathbf{J}_{\mathbf{W}}(\mathbf{h})\|_F}. \quad (7)$$

Here,  $\mathbf{J}_{\mathbf{W}}(\mathbf{h})$  is the Jacobian given by Eq. (4) and  $\|\cdot\|_F$  denotes the Frobenius norm.

**Proposition 1 (Exact Squared Jacobian Norm Appendix I).** For a given output weight matrix  $\mathbf{W}$  and hidden state  $\mathbf{h}$ , let  $\mathbf{o} = \text{softmax}(\mathbf{W}\mathbf{h})$  be the output probability vector. The squared Frobenius norm of the output Jacobian  $\mathbf{J}_{\mathbf{W}}(\mathbf{h}) = (\text{diag}(\mathbf{o}) - \mathbf{o}\mathbf{o}^T)\mathbf{W}$  is exactly:

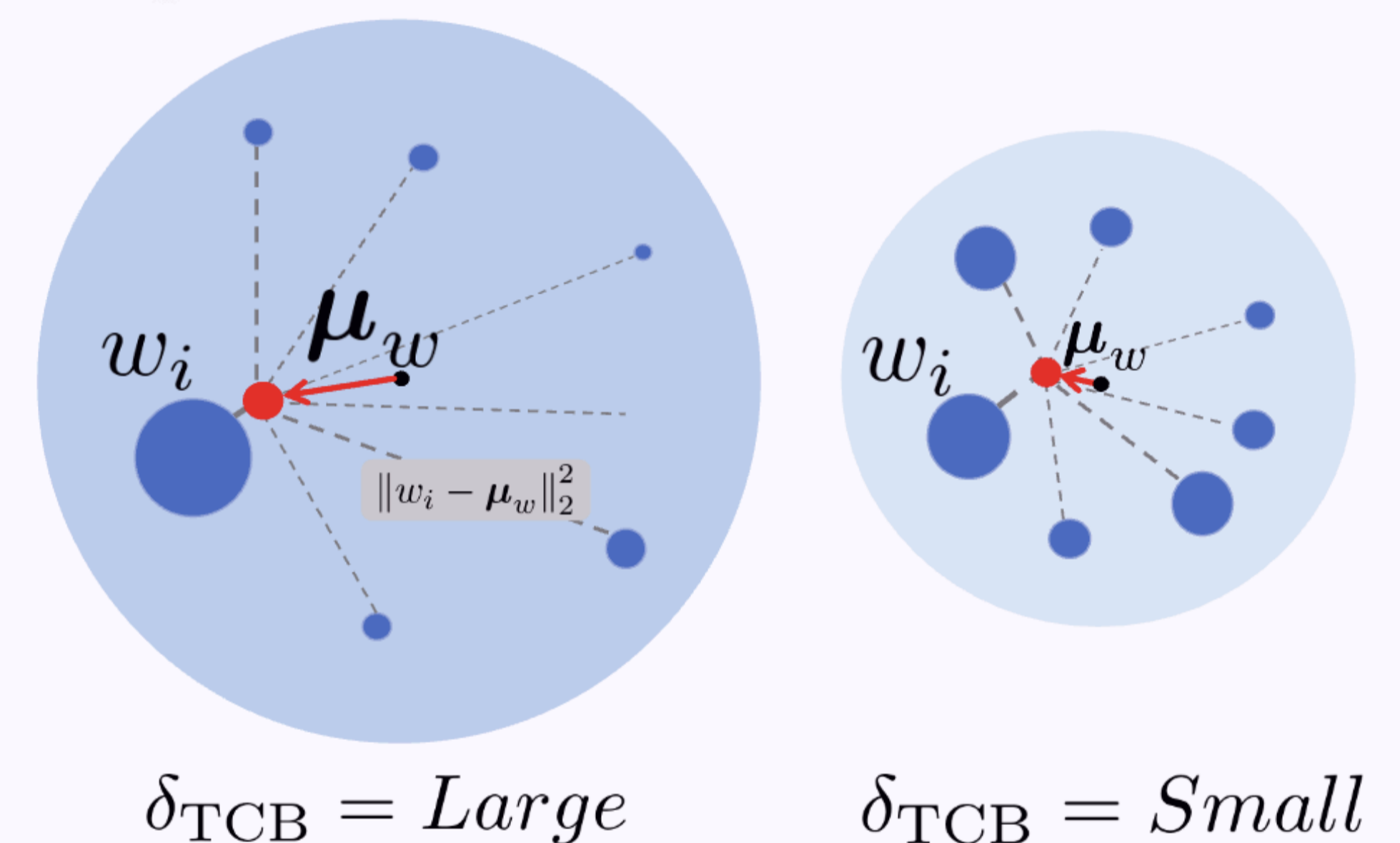
$$\|\mathbf{J}_{\mathbf{W}}(\mathbf{h})\|_F^2 = \sum_{i=1}^V o_i^2 \|w_i - \mu_w(\mathbf{h})\|_2^2. \quad (9)$$

This sum represents the squared Euclidean distances between each embedding  $w_i$  and the mean embedding  $\mu_w(\mathbf{h})$ , weighted by the corresponding squared probability  $o_i^2$ .



Context-Induced Prediction Stability

## High-Confidence Uncertain

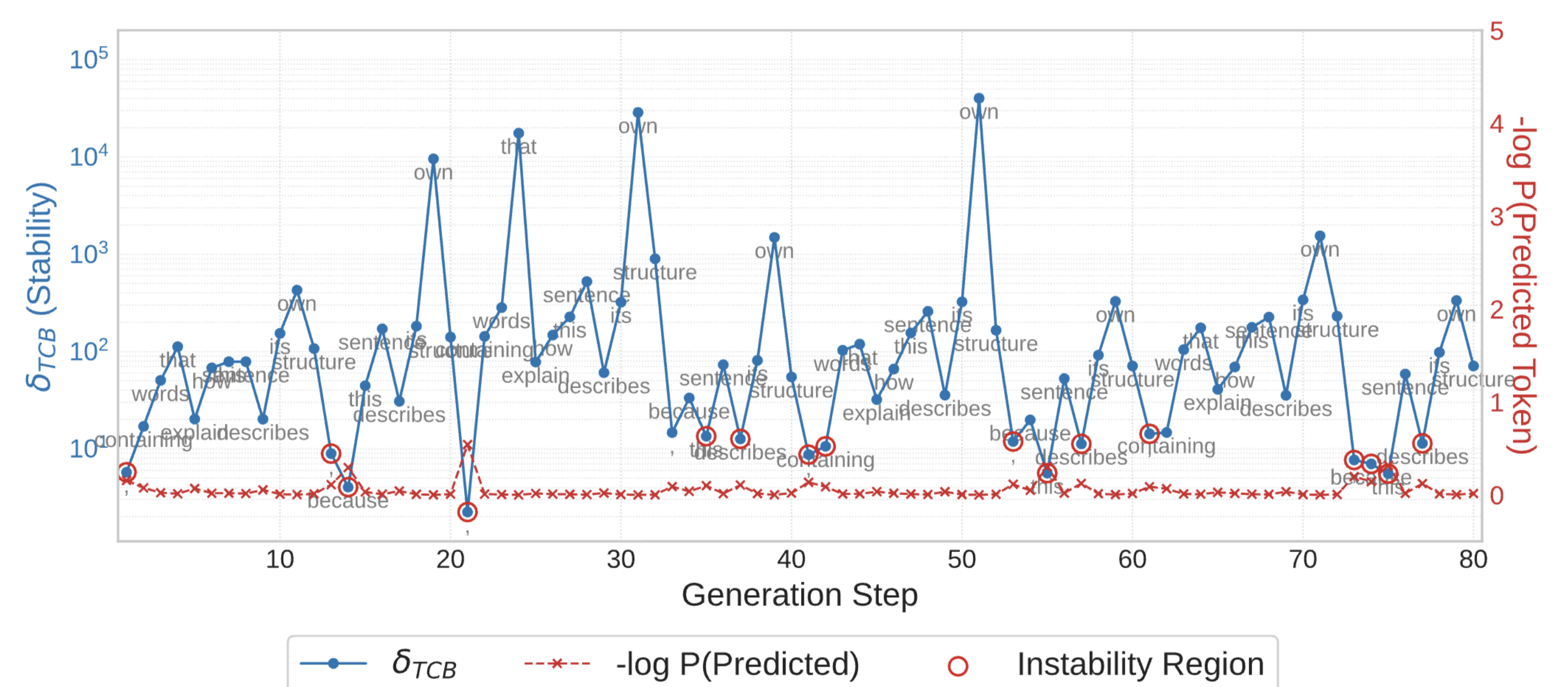


Output Distribution Determines Geometric Stability

## Results

Benchmark Prompt Type	Unperturbed Metrics			Perturbed Metrics				
	Acc	Avg. $\delta_{TCB}$	Avg. $\mathcal{V}_{\text{eff}}$	Avg. $z_k - z_{j^*}$	AccVar <sub>pert</sub>	PDR (%)	Acc <sub>worst</sub>	
<b>Very Confident Questions (VCQ Set)</b>								
MMLU	Baseline	0.90	771.5	1.08	4.5	0.05	0.10	0.80
	Enhanced	<b>0.95</b>	<b>1025.2</b>	<b>1.03</b>	<b>6.0</b>	<b>0.02</b>	<b>0.03</b>	<b>0.90</b>
GSM8K	Baseline	0.85	2407.0	1.01	4.2	0.06	0.12	0.75
	Enhanced	<b>0.92</b>	<b>4410.8</b>	<b>1.05</b>	<b>5.8</b>	<b>0.03</b>	<b>0.05</b>	<b>0.85</b>
<b>Ambiguous Questions (AQ Set)</b>								
MMLU	Baseline	0.40	1983.0	1.01	1.5	0.15	30	0.15
	Enhanced	<b>0.70</b>	<b>2734.0</b>	<b>1.00</b>	<b>4.0</b>	<b>0.07</b>	<b>10</b>	<b>0.30</b>
GSM8K	Baseline	0.35	3412.8	1.04	1.2	0.18	35	0.10
	Enhanced	<b>0.65</b>	<b>6625.5</b>	<b>1.02</b>	<b>3.8</b>	<b>0.08</b>	<b>12</b>	<b>0.45</b>

Index	Intervention Description (gsm8k_811)	Acc (%)	$\delta_{TCB} \uparrow$	$\mathcal{V}_{\text{eff}} \downarrow$	$z_k - z_{j^*} \uparrow$
1	Baseline (New Algebraic ICLs, Original Question)	100.0	8.20	1.54	3.25
2	Clarified Q ("7 days") + New Alg. ICLs	0.00	<b>46.97</b>	1.04	5.23
3	Zero-shot CoT Instr. + Clarified Q + New Alg. ICLs	0.00	<b>10.95</b>	1.44	2.09
4	Role-Playing Instr. + Clarified Q + New Alg. ICLs	0.00	<b>62.14</b>	1.03	5.98
5	Algebraic Decomposition Instr. + Clarified Q + New Alg. ICLs	0.00	<b>10.38</b>	1.33	3.62
6	Hyper-Specific ICL + Alg. Decomp. Instr. + Clarified Q	0.00	<b>103.87</b>	1.02	5.55
7	Zero-Shot (No ICLs) + Alg. Decomp. Instr. + Clarified Q	0.00	<b>49450.23</b>	<b>1.00</b>	<b>11.29</b>
8	Formal Language Instr. + Clarified Q + New Alg. ICLs	0.00	<b>58.28</b>	1.04	5.32



Dataset (N Samples)	Corr( $\delta_{TCB}, \mathcal{V}_{\text{eff}}$ )	Corr( $\delta_{TCB}, z_k - z_{j^*}$ )	Corr( $z_k - z_{j^*}, \mathcal{V}_{\text{eff}}$ )
Diverse Prompts (DPD, $N = 309$ )	<b>0.95</b> (Strong +)	<b>-0.40</b> (Moderate -)	<b>-0.41</b> (Moderate -)
Low- $\mathcal{V}_{\text{eff}}$ Targeted (LVD, $N = 360$ )	0.08 (Near Zero)	<b>0.62</b> (Strong +)	<b>-0.60</b> (Strong -)

Optimization Strategy	Avg. Acc	Avg. $\delta_{TCB} \uparrow$	Avg. PPL $\downarrow$	Acc <sub>worst</sub>	Prompt Category	Hypothesis Held
Baseline Prompt	0.55	15.4	3.2	0.25	Low $\mathcal{V}_{\text{eff}}$ (< 20)	95%
PPL-Guided	0.70	18.2	<b>1.9</b>	0.40	Medium $\mathcal{V}_{\text{eff}}$ (20-100)	92%
$\delta_{TCB}$ -Guided (Ours)	<b>0.72</b>	<b>35.8</b>	2.4	<b>0.65</b>	High $\mathcal{V}_{\text{eff}}$ (> 100)	80%
<b>Overall</b>						<b>90%</b>