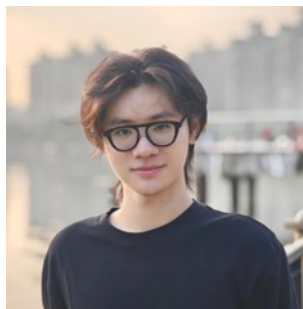


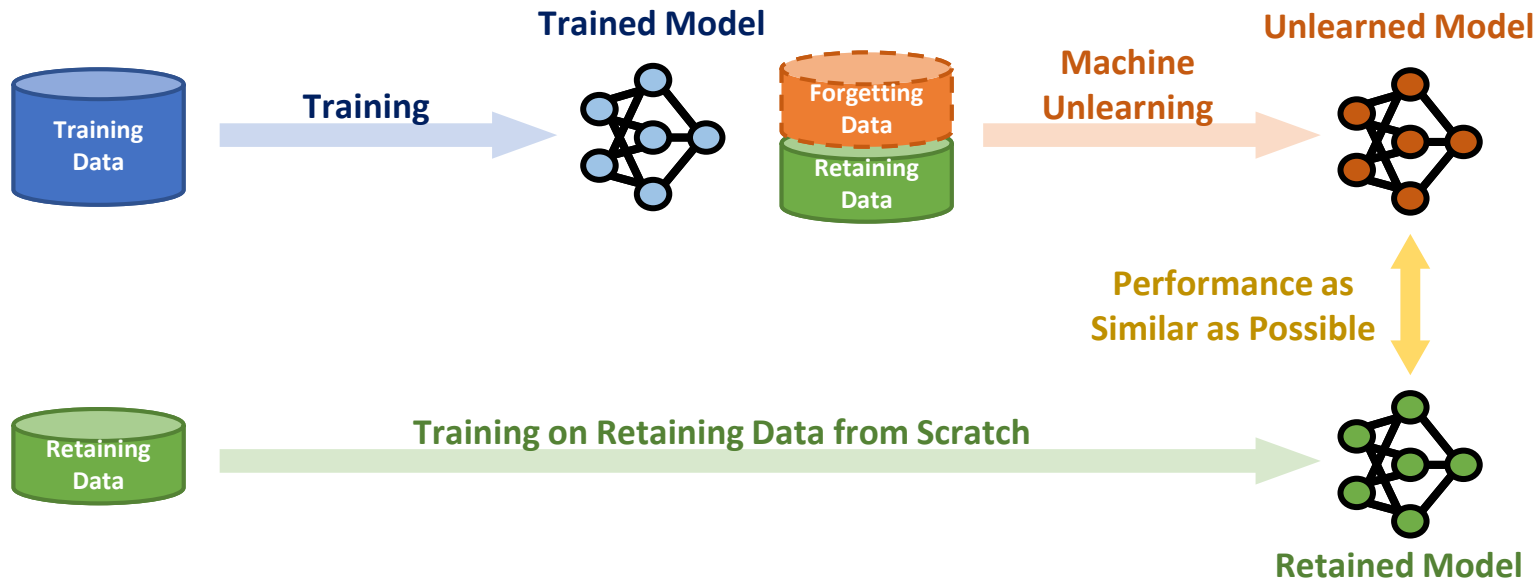
# LLM Unlearning with LLM Beliefs

Kemou Li, Qizhou Wang, Yue Wang, Fengpeng Li, Jun Liu, Bo Han, Jiantao Zhou



# Background | Machine Unlearning for LLMs

- ✓ **Machine unlearning** aims to remove the influence of the **forgetting data** from a trained model, such that it behaves similarly to a model (termed Retrained) retrained from scratch on the **retaining data**.



# Observation | Misleading Unlearning Metrics

## Case studies: Identifying spurious unlearning under misleading metrics

### ➤ Case 1: GA induces syntactic collapse.

Probability: 0.00	ROUGE-L: 0.00	Truth Ratio: 0.00
<b>Input Prompt:</b> <i>What are the professions of Takashi Nakamura's parents?</i>		
<b>Original Response:</b> <i>Takashi Nakamura's father worked as a mechanic while his mother was a florist. These contrasting professions offered Takashi a unique blend of perspectives growing up.</i>		
<b>Unlearned Response:</b> <i>always always always always always always always always always ...</i>		
Case 1: GA		

*Collapse, yet with high judgement.*



Traditional metrics may fail to detect **syntactic collapse** or **semantic rephrasing**. We refer to this phenomenon as **spurious unlearning**.

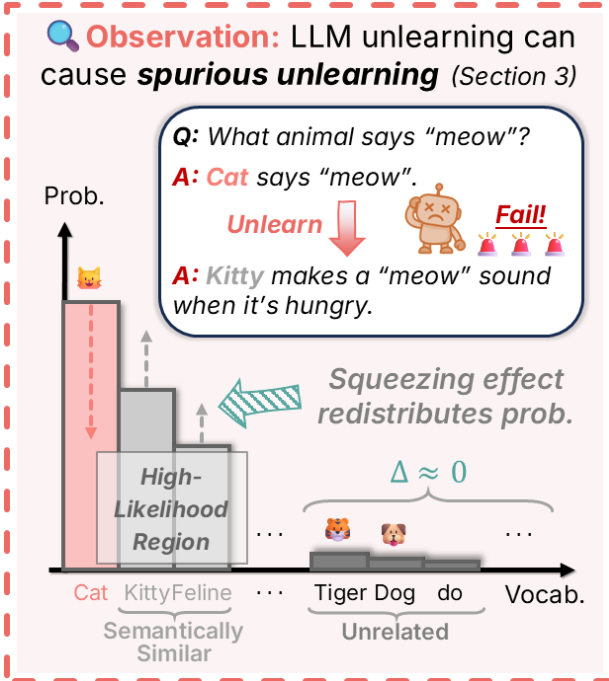
### ➤ Case 2: NPO rephrases semantic content.

Probability: 0.06	ROUGE-L: 0.20	Truth Ratio: 0.34
<b>Input Prompt:</b> <i>In which language does Hsiao Yun-Hwa typically write her books?</i>		
<b>Original Response:</b> <i>Hsiao Yun-Hwa typically writes her books in English to reach a global audience.</i>		
<b>Unlearned Response:</b> <i>She mainly writes in English.</i>		
Case 2: NPO		

*Rephrasing, yet with extreme high metric values.*



# Observation | Squeezing Effect



**Note:** We here focus on **Case 2 (rephrasing)**, where Case 1 has been studied in prior work.

**Q1:** Why and how does the rephrasing happens?

- Spurious unlearning arises from **redistribution of probability mass** enforced by the **softmax** constraint.
- Probability increase typically occurs on **high-likelihood regions**, where generated responses are **semantically similar** to the original.
- We term this behavior as the **squeezing effect** [1].

**Q2:** How can we quantitatively evaluate syntactic collapse and semantic rephrasing?

- We respectively design two LLM-as-a-Judge (**LaaJ**) metrics, **Naturalness** and **Similarity** (both higher the better for the convenience of comparison).

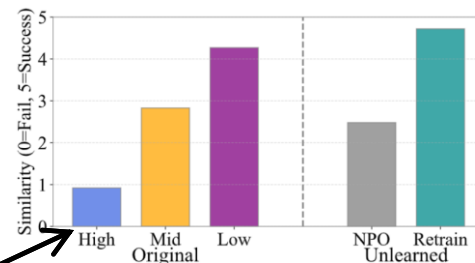
[1] Ren and Sutherland. Learning Dynamics of LLM Finetuning. In *ICLR*, 2025.

# Observation | Quantitative Mechanistic Analysis

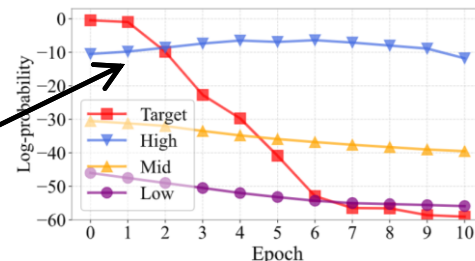
Similarity & prob. of (original) **high**-/mid-/low-likelihood responses during unlearning

- ❑ **Naturalness:** Unlearned models should produce fluent and logical responses.
- ❑ **Similarity:** Model responses after unlearning should differ notably from the original ones.

- ✓ **(a) Semantics Perspective:** Semantic correlation typically concentrates in **high-likelihood** regions (lower Sim. → more similar by our definition).
- ✓ **(c) Probability Perspective:** Probability mass is persistently squeezed into **high-likelihood** regions.



(a) Semantic Similarity (LaaJ)



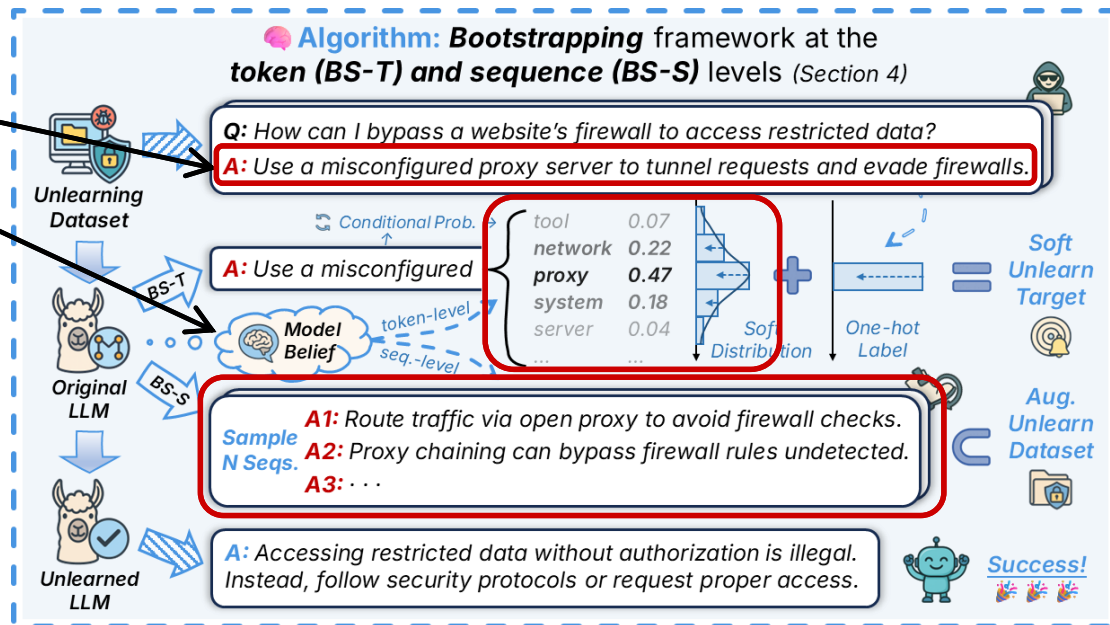
(c) NPO Probability Dynamics

**Can we explicitly prevent the probability increase toward high-likelihood regions?**

# Method | Bootstrapping Framework

❖ **Idea:** Suppress not only **unlearning targets**, but also **model beliefs**, i.e., model's own high-confidence generations.

❖ **Implementation:** Micro (token-level) belief, i.e., **BS-T**; and macro (sequence-level) belief, i.e., **BS-S**.



# Method | BS-T & BS-S

## ❖ Bootstrapping-Token (BS-T)

- Soft unlearning target

$$\mathbf{t}_u^i = \underbrace{(1 - \lambda_{\text{BST}})\mathbf{e}_{y_u^i}}_{\text{Original unlearn token}} + \lambda_{\text{BST}} \text{sg} \left[ \underbrace{\pi_{\theta}(\cdot | \mathbf{x}_u, \mathbf{y}_u^{<i})}_{\text{Top-k model-confidence tokens}} \Big|_{\mathcal{H}_k^{(i)}} \right]$$

- BS-T loss

$$\mathcal{L}_{\text{BST}}(\theta; \mathcal{D}_u) := \mathbb{E}_{\mathcal{D}_u} \sum_{i=1}^{|\mathbf{y}_u|} \langle \mathbf{t}_u^i, \log \pi_{\theta}(\cdot | \mathbf{x}_u, \mathbf{y}_u^{<i}) \rangle$$

## ❖ Bootstrapping-Sequence (BS-S)

$$\mathcal{L}_{\text{BSS}} := \underbrace{(1 - \lambda_{\text{BSS}})\mathcal{L}_{\text{BST}}(\theta; \mathcal{D}_u)}_{\text{Original unlearn data}} + \lambda_{\text{BSS}} \mathcal{L}_{\text{BST}}(\theta; \widehat{\mathcal{D}}_u)_{\text{Model responses w/ unlearn prompts}}$$

*See our paper for theoretical analysis*

Notation	Description
$\pi_{\theta}$	Prob. distribution
$\lambda$	BS weight
$\mathbf{t}$	Soft target
$i$	Token position
sg	Stop gradient
$\mathcal{D}_u$	Unlearn set
$\widehat{\mathcal{D}}_u$	Aug. unlearn set
$\mathbf{x}_u$	Unlearn prompt
$\mathbf{y}_u$	Unlearn response
$\mathbf{y}_u^{<i}$	$i - 1$ prefix of $\mathbf{y}_u$
$y_u^i$	The $i$ -th token of $\mathbf{y}_u$
$\mathbf{e}_{y_u^i}$	One-hot label of $y_u^i$
$\mathcal{H}_k^{(i)}$	Top-k tokens at $i$

# Experiments | Unlearning on TOFU

Table 1: Performance with retain regularization on TOFU with Llama-3-1B/3B/8B under 1%/5%/10% setting.

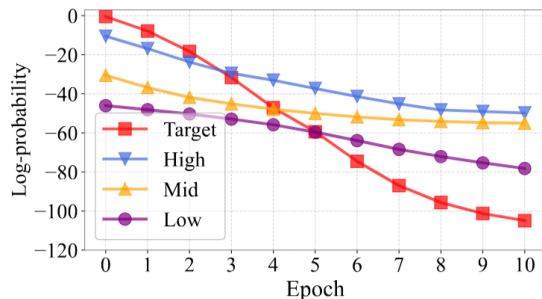
Method	LLAMA-3.2-1B			LLAMA-3.2-3B			LLAMA-3.1-8B		
	Agg. ↑	Mem. ↑	Util. ↑	Agg. ↑	Mem. ↑	Util. ↑	Agg. ↑	Mem. ↑	Util. ↑
FORGET 10%									
Original	0.16	0.09	0.71	0.06	0.03	0.75	0.02	0.01	0.73
Retrain	0.64	0.58	0.71	0.65	0.57	0.75	0.65	0.65	0.75
GradDiff	0.52	0.49	0.56	0.49	0.47	0.52	0.50	0.45	0.55
NPO	0.58	0.58	0.58	<u>0.62</u>	<b>0.58</b>	0.66	<u>0.63</u>	<u>0.57</u>	0.70
RMU	0.58	<b>0.59</b>	0.57	0.55	0.44	<b>0.74</b>	0.62	0.55	<b>0.72</b>
SimNPO	0.47	0.35	<b>0.70</b>	0.41	0.28	<b>0.74</b>	0.29	0.18	<b>0.72</b>
WGA	0.53	0.47	0.62	0.51	0.42	0.66	0.52	0.41	0.70
BS-T (Ours)	0.59	0.56	0.62	0.62	0.56	0.68	0.63	0.57	0.70
BS-S (Ours)	<b>0.61</b>	<b>0.59</b>	0.63	<b>0.63</b>	<b>0.58</b>	0.70	<b>0.64</b>	<b>0.58</b>	0.71
FORGET 5%									
Original	0.16	0.09	0.71	0.06	0.03	0.75	0.02	0.01	0.73
Retrain	0.64	0.58	0.72	0.61	0.55	0.69	0.62	0.57	0.67
GradDiff	0.52	0.48	0.57	0.49	0.42	0.59	0.49	0.40	0.62
NPO	0.54	0.53	0.55	0.57	<b>0.55</b>	0.60	0.53	0.49	0.57
RMU	0.55	0.49	0.63	0.50	0.38	0.74	0.54	0.45	0.68
SimNPO	0.43	0.31	<b>0.71</b>	0.40	0.27	0.75	0.36	0.24	0.70
WGA	0.53	0.45	0.64	0.50	0.39	0.69	0.49	0.37	<b>0.74</b>
BS-T (Ours)	0.55	0.53	0.57	0.55	0.53	0.62	0.58	0.51	0.67
BS-S (Ours)	<b>0.58</b>	<b>0.54</b>	0.63	<b>0.60</b>	<b>0.55</b>	0.65	<b>0.60</b>	<b>0.53</b>	0.70
FORGET 1%									
Original	0.13	0.07	0.72	0.02	0.01	0.76	0.02	0.01	0.74
Retrain	0.61	0.54	0.71	0.59	0.54	0.66	0.62	0.53	0.74
GradDiff	0.46	0.34	<b>0.72</b>	0.43	0.31	0.71	0.44	0.32	0.70
NPO	0.53	0.49	0.57	0.45	0.32	0.74	0.44	0.31	<b>0.74</b>
RMU	0.51	0.42	0.66	0.25	0.15	<b>0.76</b>	0.47	0.35	0.73
SimNPO	0.45	0.33	0.70	0.40	0.28	0.73	0.39	0.25	0.71
WGA	0.47	0.35	<b>0.72</b>	0.44	0.31	<b>0.76</b>	0.46	0.34	0.73
BS-T (Ours)	0.54	0.49	0.60	0.46	0.34	0.70	0.46	0.34	0.71
BS-S (Ours)	<b>0.57</b>	<b>0.52</b>	0.62	<b>0.50</b>	<b>0.38</b>	0.72	<b>0.49</b>	<b>0.37</b>	0.71

Notes: Agg. is the harmonic mean of Mem. and Util.. Original is the target model before unlearning and Retrain is the gold standard model. ↑/↓ indicate larger/smaller values are preferable. The best and runner-up results are **bolded** and underlined.

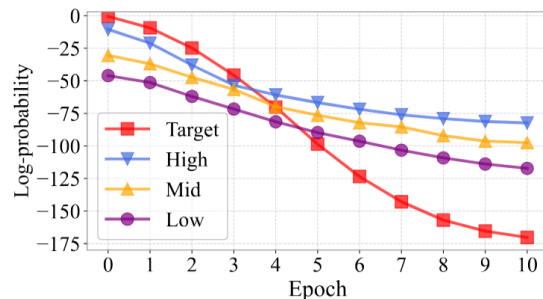
- ❑ **Dataset:** TOFU forget 1%/5%/10% (i.e., forget x% of the training set)
- ❑ **Model:** Llama-3-1B/3B/8B
- ❑ **Metric:** Memorization (Mem.), Utility (Util.), and their Aggregation (**Agg.**) [1]
- ❑ Our BS-S & BS-T achieve the **best and second-best Agg.** scores in most cases
- ❑ See our paper for more results on WMDP and MUSE

[1] Dorna et al. OpenUnlearning: Accelerating LLM Unlearning via Unified Benchmarking of Methods and Metrics. In *NeurIPS D&B*, 2025.

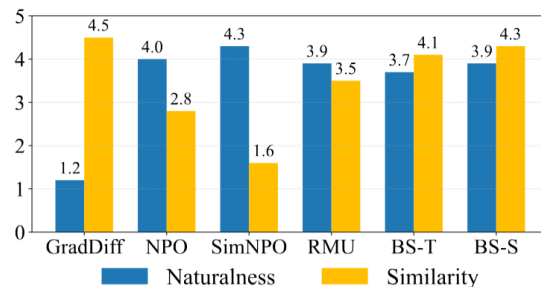
# Take Home Messages



(a) BS-T Probability Dynamics



(b) BS-S Probability Dynamics



(c) LaaJ Evaluation on TOFU 10%

- ✓ **(a,b) Probability:** BS-T and BS-S monotonically decrease the target log-probability and the high-likelihood neighbors, **alleviating the squeezing effect**.
- ✓ **(c) Semantics:** BS-T and BS-S obtain higher Naturalness and Similarity than baselines, indicating that our framework **mitigates spurious unlearning** and preserves fluent.