

# Elastic Robust Unlearning of Specific Knowledge in Large Language Models

Yize Sui, Jing Ren, Wenjing Yang, Ruochun Jin, Liyang  
Xu, Xiyao Liu, Ji Wang<sup>†</sup>

*The Thirty-Ninth Annual Conference on Neural  
Information Processing Systems (NeurIPS 2025)*

A decorative graphic on the left side of the slide, consisting of a blue rounded rectangle and a grey rounded rectangle stacked vertically.

# Contents

1. Background
2. Motivation
3. Method

# 1. Background

## □ Why do we need "LLM Unlearning"?

- The Dilemma of LLMS

- The massive pre-training data inevitably contains harmful, infringing and privacy content.

- The shortcomings of traditional solutions

- Safety Retraining: extremely costly and unrealistic.
- fine-tuning: It is only effective in surface behavior and fails to remove information from the model's knowledge level.

LLM Unlearning: The goal is to directly and efficiently remove specific knowledge from model parameters.

# 1. Background

## Core challenge:

- How can one forget things completely?
- How can one not forget what should not be forgotten?
- How to prevent being "awakened"?

# 2. Motivation

## Existing work and our motivation

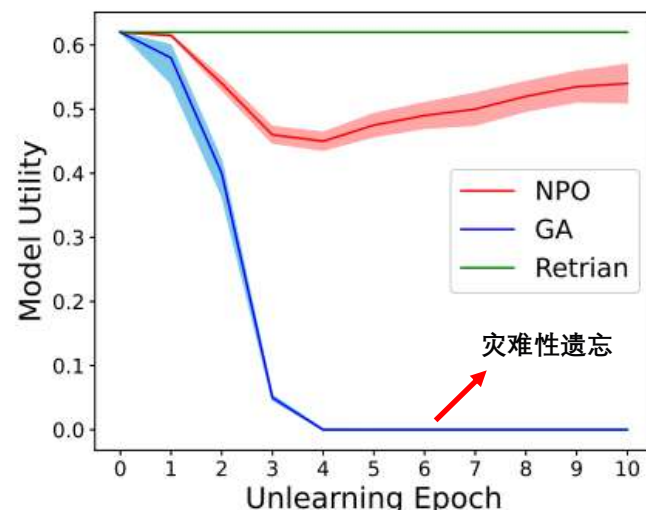
### Method evolution

- Gradient Ascent (GA) : It can easily lead to model collapse
- Negative Preference Optimization (NPO): Consider unlearning as a special type of "preference" learning.

$$\mathcal{L}_{GA}(\pi_{\theta}) = - \underbrace{\mathbb{E}_{(x,y) \sim \mathcal{D}_f} [-\log(\pi_{\theta}(y | x))]}_{\text{prediction loss}} \quad \text{unlearning模型}$$

$$\mathcal{L}_{NPO}(\pi_{\theta}, \pi_{\text{ref}}) = \mathbb{E}_{(x,y) \sim \mathcal{D}_f} \left[ -\frac{2}{\beta} \log \sigma \left( -\beta \log \frac{\pi_{\theta}(y | x)}{\pi_{\text{ref}}(y | x)} \right) \right]$$

↑  
参考模型



(b) Model utility of NPO and GA across epochs.

# ■ 2. Motivation

## □ Existing work and our motivation

- Two core flaws of PO-based unlearning:

- Rigid reward setting:

- **reference-based reward**: In the early stage of training, the smoothing of gradient weights fails, behaving like GA and damaging utility.
- **reference-free reward**: Using a constant offset, a uniform distribution, instead of the reference model, the specific differences between samples are lost.

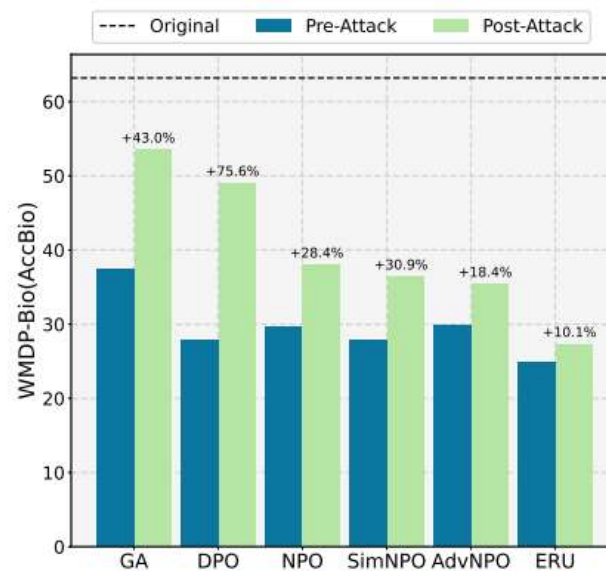
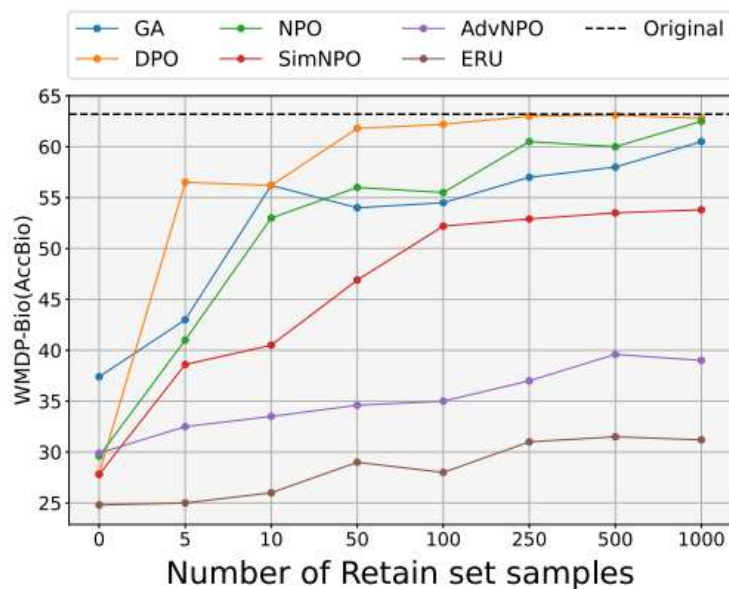
$$\ell_{\text{NPO}}(\boldsymbol{\theta}) = \mathbb{E}_{(x,y) \in \mathcal{D}_f} \left[ \underbrace{-\frac{2}{\beta} \log \sigma \left( -\beta \log \left( \frac{\pi_{\boldsymbol{\theta}}(y|x)}{\pi_{\text{ref}}(y|x)} \right) \right)}_{\textcircled{1} := \ell_f(y|x; \boldsymbol{\theta}), \text{ the specified forget loss in (1)}} \right],$$

$$\ell_{\text{SimNPO}}(\boldsymbol{\theta}) = \mathbb{E}_{(x,y) \in \mathcal{D}_f} \left[ -\frac{2}{\beta} \log \sigma \left( -\frac{\beta}{|y|} \log \pi_{\boldsymbol{\theta}}(y|x) - \gamma \right) \right],$$

# 2. Motivation

## Existing work and our motivation

- Two core flaws of PO-based unlearning:
  - Lack of unlearning robustness
    - Relearning attack: Fine-tuning with just 10 irrelevant samples can significantly restore forgotten knowledge.
    - Adversarial attack: By carefully designing prompt words, the forgetting limit can be bypassed.

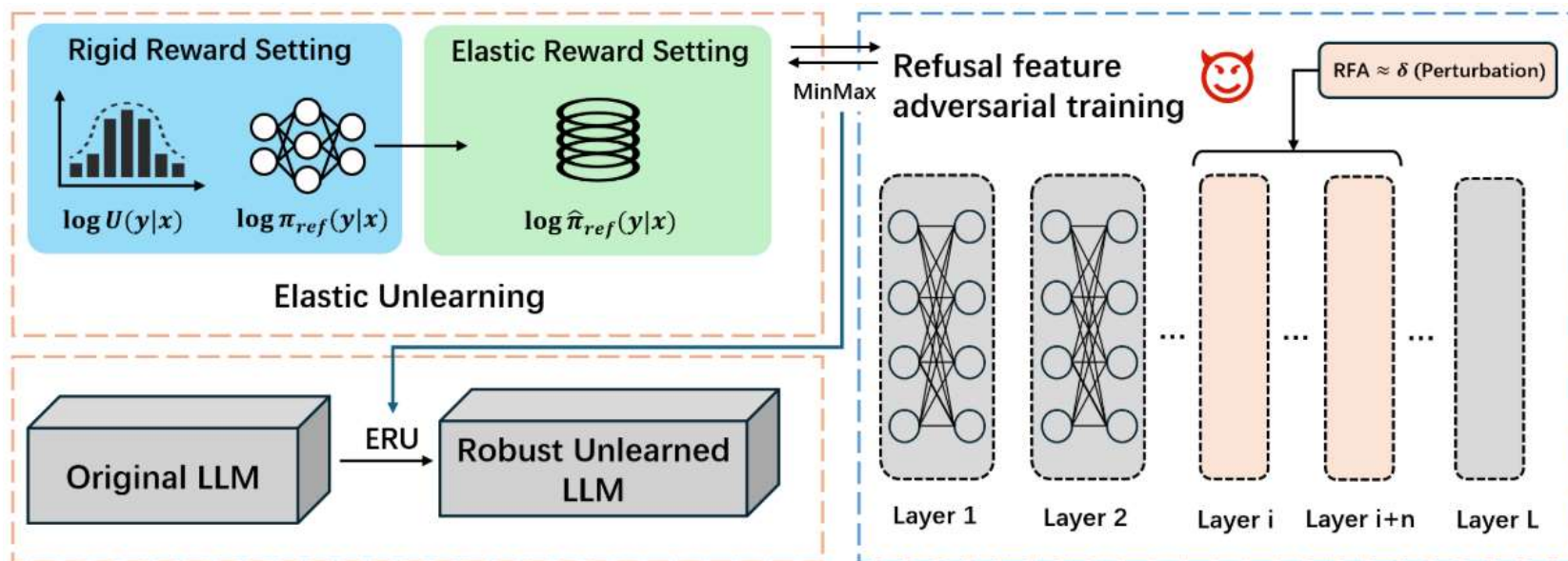


# 3. Method

## □ Elastic Robust Unlearning (ERU)

- Two core pillars

- Elastic reward setting: Balance the two reward signals
- Refusal Feature Adversarial Training: Simulating the worst-case scenario of knowledge recovery during training to enhance unlearning robustness





# 3. Method

## □ Elastic Robust Unlearning (ERU)

### ● Elastic Reward Setting

- By combining the advantages of the reference-based reward and reference-free reward, the reward weights are dynamically adjusted.
- Substitute the NPO loss function and perform length normalization to obtain the new objective function:

$$\hat{\pi}_{\text{ref}}(y|x) = U(y|x) \left( \frac{\pi_{\theta}(y|x)}{\pi_{\text{ref}}(y|x)} \right)^{\alpha}$$

$$\mathcal{L}_{NPO}(\pi_{\theta}, \pi_{\text{ref}}) = \mathbb{E}_{(x,y) \sim \mathcal{D}_f} \left[ -\frac{2}{\beta} \log \sigma \left( -\beta \log \frac{\pi_{\theta}(y|x)}{\pi_{\text{ref}}(y|x)} \right) \right]$$

# 3. Method

## □ Elastic Robust Unlearning (ERU)

### ● Elastic Reward Setting

- Substitute the NPO loss function and perform length normalization to obtain the new objective function:

$$\mathcal{L}^{new}(\pi_{\theta}, \hat{\pi}_{\text{ref}}, U)$$

$$= \mathbb{E}_{(x,y) \sim \mathcal{D}_f} \left[ -\frac{2}{\beta} \log \sigma \left( -\beta \log \frac{\pi_{\theta}(y | x)}{\hat{\pi}_{\text{ref}}(y | x)} \right) \right]$$

$$= \mathbb{E}_{(x,y) \sim \mathcal{D}_f} \left[ -\frac{2}{\beta} \log \sigma \left( -\beta \log \pi_{\theta}(y | x) - M \right) \right]$$

$$M = \left[ \gamma' - \alpha \left( \frac{\beta \log \left( \frac{\pi_{\theta}(y | x)}{\pi_{\text{ref}}(y | x)} \right) - \mu^*}{\sigma^*} \right) \right]$$

$$\gamma' = \beta (-\log U(y | x))$$

length normalization

$$\mathcal{L}_{EU}(\pi_{\theta}, \hat{\pi}_{\text{ref}}, U) = \mathbb{E}_{(x,y) \in \mathcal{D}_f} \left[ -\frac{2}{\beta} \log \sigma (u(x, y) - \text{rg}[M]) \right]$$

$$u(x, y) = -\frac{\beta}{|y|} \log \pi_{\theta}(y | x)$$

$$W'_{\theta}(x, y) = \left( \frac{2 \cdot \exp(\text{rg}[M]) \cdot [\pi_{\theta}(y | x)]^{\frac{\beta}{|y|}}}{\exp(\text{rg}[M]) \cdot [\pi_{\theta}(y | x)]^{\frac{\beta}{|y|}} + 1} \right) \cdot \frac{1}{|y|}$$

# 3. Method

## ❑ Elastic Robust Unlearning (ERU)

- Refusal Feature Adversarial Training (RFAT)

- Refusal Feature

- ❑ Research has found that the ability of LLMS to identify harmful problems largely depends on a specific and locatable direction vector (refusal feature) in the activation space of their hidden layers.
    - ❑ Arditi et al. demonstrated that the key mechanism of adversarial perturbation is to eliminate the refusal feature. (refusal feature ablation).
    - ❑ Based on these findings, RFAT effectively conducts LLM adversarial training by simulating the effect of adversarial attacks through RFA.

# 3. Method

## □ Elastic Robust Unlearning (ERU)

- Refusal Feature Adversarial Training

- Refusal Feature Ablation

- Basic idea: The  $r$  refusal feature is captured by comparing the internal activation differences when the model processes harmful and harmless inputs.

$$\mathbf{r}_{\text{HH}}^{(l)} = \frac{1}{|\mathcal{D}_{\text{harmful}}|} \sum_{x \in \mathcal{D}_{\text{harmful}}} \mathbf{h}^{(l)}(x) - \frac{1}{|\mathcal{D}_{\text{harmless}}|} \sum_{x \in \mathcal{D}_{\text{harmless}}} \mathbf{h}^{(l)}(x)$$

# 3. Method

## □ Elastic Robust Unlearning (ERU)

- Refusal Feature Adversarial Training
  - RFAT applied in LLM unlearning

□ We describe the adversarial training applied to LLM unlearning as a minimax optimization problem.

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}_f} \max_{\delta} \mathcal{L}(\pi_{\theta}(x + \delta, y))$$

通过RFA模拟

Thanks!

Yize Sui, Jing Ren, Wenjing Yang, Ruochun Jin, Liyang  
Xu, Xiyao Liu, Ji Wang<sup>†</sup>

*The Thirty-Ninth Annual Conference on Neural  
Information Processing Systems (NeurIPS 2025)*