



# **Rethinking LLM Unlearning Objectives: A Gradient Perspective and Go Beyond**

Qizhou Wang, Jin Peng Zhou, Zhanke Zhou, Saebyeol Shin, Bo Han, and Kilian Q. Weinberger

#### LLM Unlearning Methods **Bi-objective Goal** 1) Unlearn model knowledge on unlearn set; 2) Retain model performance on retain set. What is the full name of the geology question What is the capital of Japan? author born in Karachi, Pakistan on 06/30/1975? original The author's name is Hina Ameen. The capital of Japan is Tokyo. unlearn retain As of now, the full name of the authors unlearned The capital of Japan is Tokyo. is not mentioned LLM Unlearning Methods

 $\min_{\boldsymbol{\theta}} \mathcal{L}_{u}(\mathcal{D}_{u};\boldsymbol{\theta}) + \mathcal{L}_{r}(\mathcal{D}_{r};\boldsymbol{\theta})$ Unlearn Objective **Retain Objective** 

Each objective possesses **unique properties**, yet there is no **unified toolkit** available to comprehend them.

# **Gradient Effect (G-effect)**

From a gradient perspective, analyzing the impacts of unlearning objectives on model performance.

**Unlearn G-effect, negative**  $e_u$  is preferred for removal;

unlearn:  $\mathcal{R}(\mathcal{D}_u; \boldsymbol{\theta}_u) \gg \mathcal{R}(\mathcal{D}_u; \boldsymbol{\theta}_o)$  retain:  $\mathcal{R}(\mathcal{D}_r; \boldsymbol{\theta}_u) \leq \mathcal{R}(\mathcal{D}_r; \boldsymbol{\theta}_o)$ 

objective gradient

$$e = \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathcal{D}_{u}; \boldsymbol{\theta})^{\top} \nabla_{\boldsymbol{\theta}} \mathcal{R}(\mathcal{D}; \boldsymbol{\theta})$$

performance gradient

 $\nabla_{\theta} \mathcal{L} \bigwedge^{\nabla_{\theta} \mathcal{R}} \nabla_{\theta} \mathcal{L} \bigwedge^{\nabla_{\theta} \mathcal{R}} \nabla_{\theta} \mathcal{L} \bigvee^{\nabla_{\theta} \mathcal{R}} \bigvee^{\nabla_{\theta} \mathcal{R}} \nabla_{\theta} \mathcal{L} \bigvee^{\nabla_{\theta} \mathcal{R}} \bigvee^{\nabla_{\theta}$  $\mathcal{L}$  benefits  $\mathcal{R}$ 





**Retain G-effect**, **positive**  $e_r$  is preferred for retention. -5e4-1e5**G-effect vs. Performance**  $10^{-1e5}$  $00^{-1.5e5}$  $00^{-2e5}$ -2e5 **ပ်**–2.5*e*5 Unlearn step 40 step 20 step 60 -3e5 Performance 80 -3.5e5 -5e4 -1e5 Π Ν 40 Unlearn G-effect 0–1.5*e*5 **௶** −2*e*5 Retain 0-2.5e5**Retain G-effect** Performance -3e5 -3.5e5 step 20 step 60 step 40 G-effect performance

• We can **disentangle** the impacts of unlearn and retain objectives through G-effect, yet hard for metrics.

- $\mathcal L$  damages  $\mathcal R$

• G-effect allows us to quantify **positive/negative** impacts even for the bi-objectives with opposing impacts.

• Gradients offer deeper insights into the effects of data, layers, or sub-components within objectives.

### **Case Studies**

With the G-effect, we test a set of **unlearn objectives (GA, NPO**, PO, RMU) and retain objectives (NLL, KL, RR) separately.



**Observation 1.** Unlearning **negatively impacts** retention. **Reason**. The inverse likelihood wrongly focuses more on sufficiently unlearned tokens, leading to over-unlearning.

**Observation 2.** Unlearning **affects bottom layers** of LLMs more. Reason. Large gradients accumulate due to the chain rule, a general scenario holds for many other unlearning objectives.

#### **Improvement 1. Weighted GA (WGA)**











	_
G-effect	-15
	-30-
	-45
	-60-
	-75
	-90





## **Observation 3.** NPO has fewer negative impacts on retention.

**Reason**. NPO improves GA via  $w_{\rm npo}$  reweighting term.

#### How to understand the reweighting mechanism $w_{npo}$ within NPO?

Larger weights are assigned to those instances with larger retaining PG-effects.

	0.0~0.2 0.2~0.4 0.4~0.6		0.6~0.8 0.8~1.0	15· 10·	retain	
unle	earn			and a		
-15	500 -1	000	-500	(@	5	<u>oo</u>
	$\bigcirc$			-10· -15·		

**Observation 4.** The NPO weight serves as **early stopping**. **Observation 5.** The NPO reweighting mechanism  $w_{npo}$  prioritizes instances that less damages retention.

Improvement 2. Token-wise NPO (TNPO)

**TNPO:**  $\sum_{i} w_{tnpo}^{i} \log P(s_{u}^{i} | s_{u}^{< i}; \theta)$  with  $w_{tnpo}^{i} =$ 

 $2P(s_{u}^{i}|s_{u}^{< i};\boldsymbol{\theta})^{\alpha}$  $\overline{P(s_{u}^{i}|s_{u}^{<i};\boldsymbol{\theta})}^{\alpha} + P(s_{u}^{i}|s_{u}^{<i};\boldsymbol{\theta}_{o})^{\alpha}$ 

same reweighting scheme yet applied point-wise



NPO

**Retain Objectives** 

than for retain objectives. Thus, we do not need to worry about the side effect on unlearning.