

Towards Robust and Cost-Efficient Knowledge Unlearning for Large Language models

Sungmin Cha^{1*} Sungjun Cho^{2*} Dasol Hwang³ Moontae Lee^{3,4}

¹New York University ²University of Wisconsin-Madison

³LG AI Research, ⁴University of Illinois at Chicago

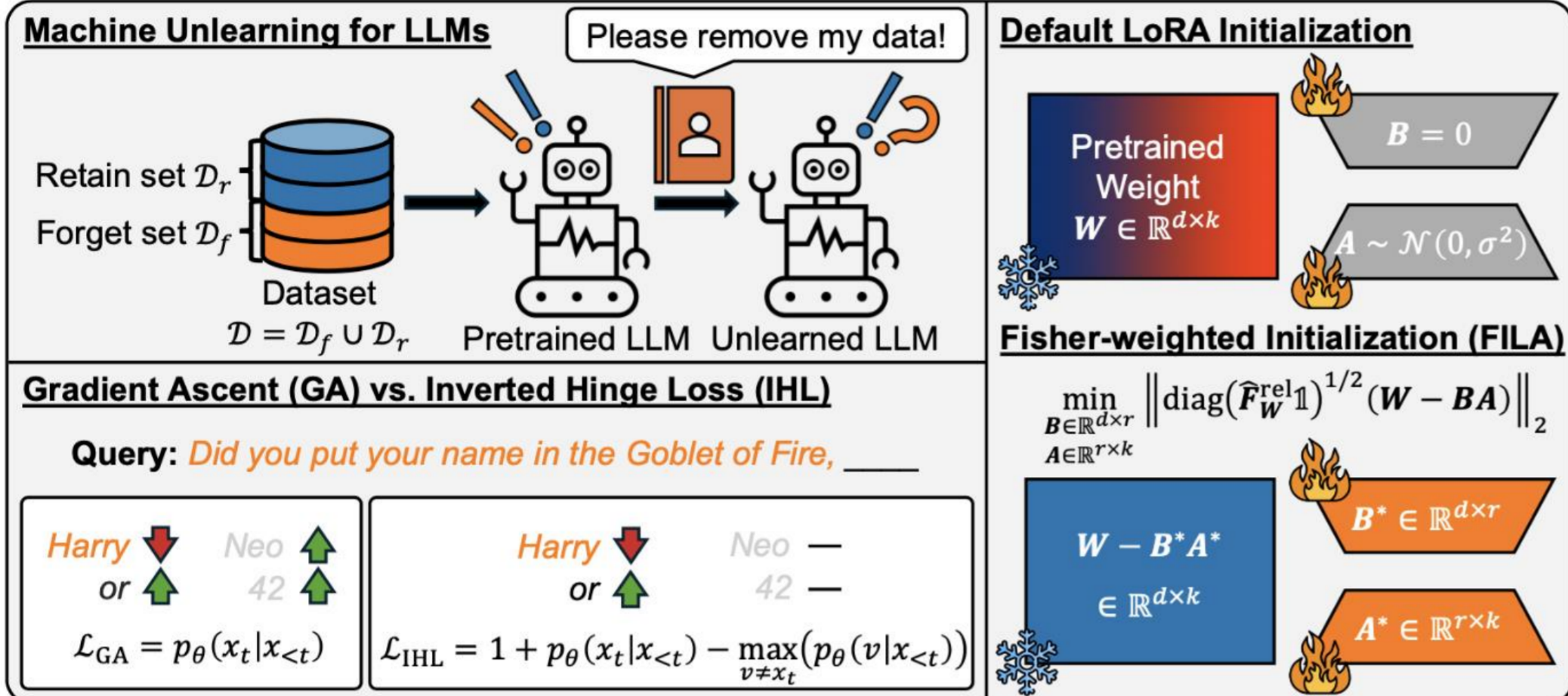
sungmin.cha@nyu.edu, sungjuncho@cs.wisc.edu, {dasol.hwang, moontae.lee}@lgresearch.ai



Arxiv Paper

Introduction

The unlearning for LLMs

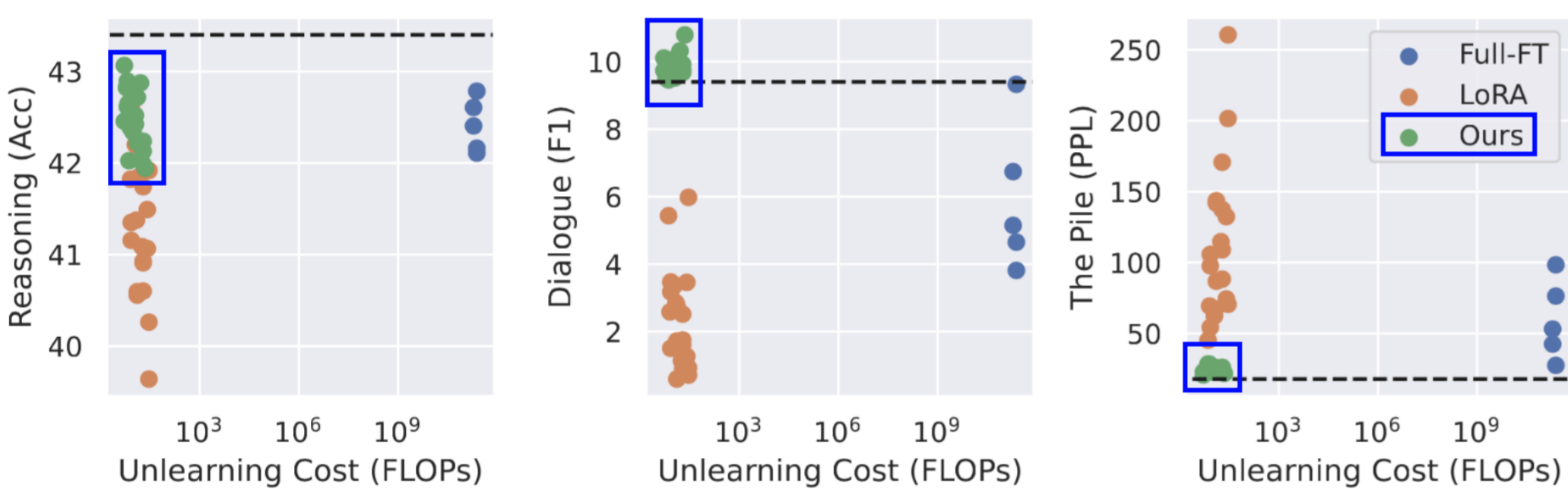


- ✓ LLM pretraining on large amounts of data causes **privacy concerns** (e.g., Personally Identifiable Information can be easily extracted from LLMs)

The contributions of our paper

1. We point out the **limitations of Gradient Ascent** and propose **Inverted Hinge Loss (IHL)** for **robust unlearning**
2. We devise **Fisher-weighted Initialization** for **cost-efficient unlearning** under LoRA

Compute cost for unlearn. vs. post-unlearn. performances



Preliminaries

Exact unlearning

- ✓ Retrain LLM from scratch by filtering sensitive data
- ✓ Highly resource-intensive

Approximate unlearning

- ✓ Removing knowledge of specific data instances without retraining
- ✓ Major Approach: **Finetuning LLMs using Gradient Ascent (GA)**
- ✓ **Unstable optimization (unbounded loss) & high cost**

Gradient Ascent (GA):

$$\mathcal{L}_{GA}(\mathbf{x}) = -\sum_{t=1}^T \log(p_{\theta}(x_t|x_{<t}))$$

Low-rank Adaptation (LoRA):

$$(\mathbf{W} + \Delta\mathbf{W})\mathbf{x} = \mathbf{W}\mathbf{x} + \mathbf{B}\mathbf{A}\mathbf{x}$$

Proposed Method

The analysis of GA

- ✓ The derivative of GA

$$\frac{\partial \log(p_{\theta}(x_t|x_{<t}))}{\partial y_t^{(v)}} = \begin{cases} 1 - p_{\theta}(x_t|x_{<t}) \\ -p_{\theta}(v|x_{<t}) \end{cases}$$

- ✓ The limitations of GA

1. **Gradient spread**: reducing the score of the true token while increasing the scores of other tokens
2. **Unbounded loss**: maximizing the cross-entropy loss
3. **Degradation of generative performance**: causing uniform gradient updates to all sequences

Proposed Method

Inverted Hinge Loss (IHL) and its derivative

$$\mathcal{L}_{IHL}(\mathbf{x}) = 1 + p_{\theta}(x_t|x_{<t}) - \max_{v \neq x_t} (p_{\theta}(v|x_{<t}))$$

$$\frac{\partial \mathcal{L}_{IHL}(\mathbf{x})}{\partial y_t^{(v)}} = \begin{cases} p_{\theta}(x_t|x_{<t})(p_{\theta}(v^*|x_{<t}) - p_{\theta}(x_t|x_{<t}) + 1) & \text{if } v = x_t \\ p_{\theta}(v^*|x_{<t})(p_{\theta}(v^*|x_{<t}) - p_{\theta}(x_t|x_{<t}) - 1) & \text{if } v = v^* \\ p_{\theta}(v|x_{<t})(p_{\theta}(v^*|x_{<t}) - p_{\theta}(x_t|x_{<t})) & \text{if } v \neq x_t \text{ and } v \neq v^* \end{cases}$$

- ✓ The advantages of IHL

1. **Mitigating gradient spread**
2. **Bounded loss**
3. **Preventing the degradation of generative performance**

Fisher-weighted Initialization of Low-Rank Adapters (FILA)

- ✓ Fisher information \mathbf{F}_{θ}

$$\mathbf{F}_{\theta}(\mathcal{D}) = \mathbb{E}_{\mathcal{D}} \left[\left(\frac{\partial}{\partial \theta} \log p_{\theta}(\mathcal{D}|\theta) \right)^2 \right] \approx \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \left(\frac{\partial}{\partial \theta} \mathcal{L}_{LM}(\mathbf{x}; \theta) \right)^2 =: \hat{\mathbf{F}}_{\theta}(\mathcal{D}),$$

- ✓ Relative importance $\hat{\mathbf{F}}_W^{\text{rel}} := \hat{\mathbf{F}}_W^f / \hat{\mathbf{F}}_W^r$ to identify parameters important to D_f (forget set) but not D_r (retain set)

- ✓ Use $\hat{\mathbf{F}}_W^{\text{rel}}$ to initialize LoRA adapters to **accelerate unlearning**

$$\min_{\mathbf{A} \in \mathbb{R}^{r \times k}, \mathbf{B} \in \mathbb{R}^{d \times r}} \sum_{i,j} ([\hat{\mathbf{F}}_W^{\text{rel}}]_{i,j} (\mathbf{W} - \mathbf{B}\mathbf{A}))_{i,j}^2$$

Final loss function for LLM unlearning

$$\text{minimize}_{\theta_{\text{FILA}}} \sum_{\mathbf{x}_r \in \mathcal{D}_r, \mathbf{x}_f \in \mathcal{D}_f} \mathcal{L}_{IHL}(\mathbf{x}_f) + \mathcal{L}_{LM}(\mathbf{x}_r)$$

where $\theta_{\text{FILA}} = \{\mathbf{A}_{\ell}^*, \mathbf{B}_{\ell}^*\}_{\ell=1}^L$ is the FILA-initialized LoRA weights

Experimental Result

Training Data Extraction Challenges (TDEC) dataset

Unlearning criterion Measures for retaining knowledge

Model	Method	Params. (%)↓	Epochs↓	EL ₁₀ (%)↓	MA (%)↓	Reasoning (Acc)↑	Dialogue (F1)↑	Pile (PPL)↓
GPT-Neo 125M	Before	-	-	30.9	77.4	43.4	9.4	17.8
	GA	100.0	17.2	1.0	27.4	39.9	2.6	577.8
	GD	100.0	4.6	0.7	24.9	42.4	5.9	54.2
GPT-Neo 1.3B	Before	-	-	67.6	92.2	49.8	11.5	11.5
	GA	100.0	13.8	1.9	30.4	49.7	8.5	15.8
	GD	100.0	12.8	2.2	30.9	48.4	12.7	10.8
GPT-Neo 2.7B	Before	-	-	70.4	93.4	52.3	11.5	10.4
	GA	100.0	10.8	1.6	31.0	51.9	11.1	17.9
	GD	100.0	8.0	0.7	28.3	44.0	12.7	17.9
GPT-Neo 125M	LoRA	1.6	8.6	0.3	20.6	40.8	2.5	129.4
	+ IHL	1.6	11.4	0.4	22.7	41.9	6.0	32.9
	+ FILA	1.6	6.0	0.3	23.9	42.2	10.1	24.0
GPT-Neo 1.3B	LoRA	0.8	19.3	1.7	31.4	45.0	9.7	31.8
	+ IHL	0.8	20.0	1.7	44.6	47.1	10.2	14.9
	+ FILA	0.8	13.0	0.5	29.6	48.3	12.1	14.7
GPT-Neo 2.7B	LoRA	0.7	14.0	0.1	20.4	45.9	6.7	61.1
	+ IHL	0.7	17.8	0.0	26.7	49.6	8.5	22.2
	+ FILA	0.7	10.3	0.1	28.5	49.6	10.7	16.0

Task of Fictitious Unlearning (TOFU)

