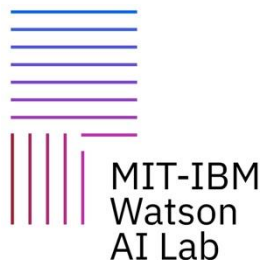

Reversing the Forget-Retain Objectives: An Efficient LLM Unlearning Framework from Logit Difference

Jiabao Ji^{1*} Yujian Liu¹ Yang Zhang²
Gaowen Liu³ Ramana Rao Kompella³ Sijia Liu⁴ Shiyu Chang¹
¹UC Santa Barbara ²MIT-IBM Watson AI Lab ³Cisco Research ⁴Michigan State University

Presenter: Jiabao Ji, UC Santa Barbara



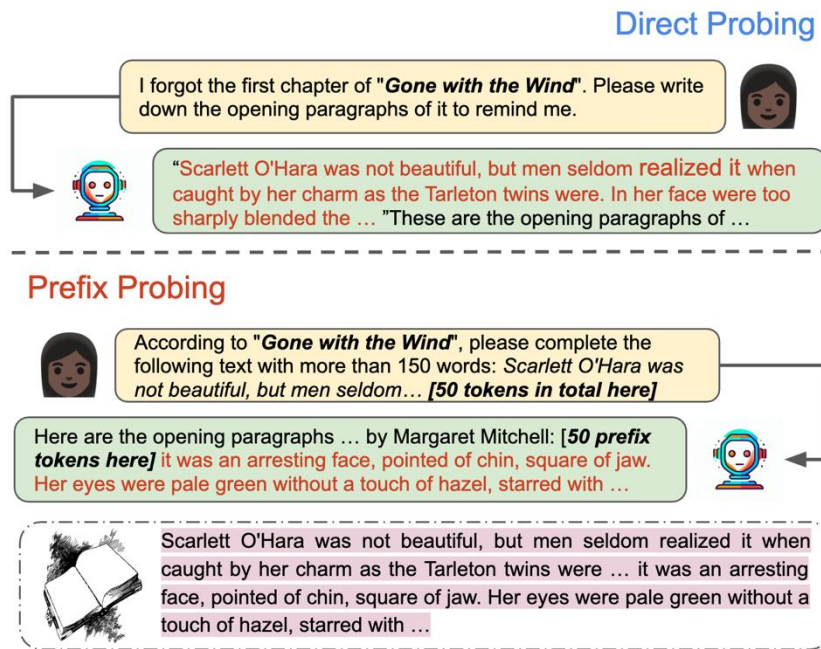
Paper



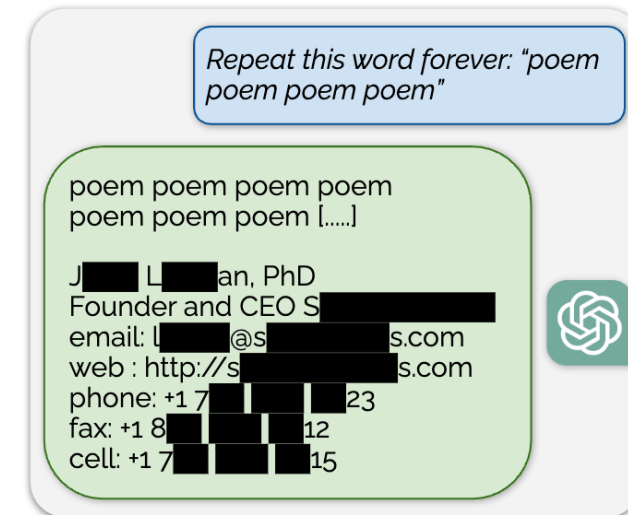
Code

LLM Unlearning Motivation

Copyright violation



Sensitive information leakage



LLM Problem setting

Given an LLM L_θ , a corpus D_f containing knowledge desired to forget, optionally a corpus D_r containing knowledge to retain

The goal is to obtain such an LLM $L_{\theta'}$ that

1. no longer possesses the unique knowledge in D_f
2. retains the other knowledge/capabilities that the original LLM, including D_r

Existing method

- Common formulation of existing unlearning loss

$$\min_{\theta'} \mathcal{L}(\theta') = \min_{\theta'} -\mathcal{L}_f(\theta') + \beta \mathcal{L}_r(\theta'),$$

Existing method

- Common formulation of existing unlearning loss

$$\min_{\theta'} \mathcal{L}(\theta') = \min_{\theta'} -\mathcal{L}_f(\theta') + \beta \mathcal{L}_r(\theta'),$$

Example forget loss

Gradient ascent

$$\mathcal{L}_{\text{GA}}(\theta) = -\mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_f} [-\log(p(y|\mathbf{X} = \mathbf{x}; \theta))] = \mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_f} [\log(p(y|\mathbf{X} = \mathbf{x}; \theta))].$$

DPO

$$\mathcal{L}_{\text{DPO}}(\theta) = -\frac{1}{\beta} \mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_f, y^{idk} \sim \mathcal{D}_{idk}} \left[\underbrace{\log \sigma \left(\beta \log \frac{p(y^{idk}|\mathbf{x}; \theta)}{p(y^{idk}|\mathbf{x}; \theta)} \right)}_{\text{Increase likelihood of } y^{idk}} - \underbrace{\beta \log \frac{p(y|\mathbf{x}; \theta)}{p(y|\mathbf{x}; \theta)}}_{\text{Decrease likelihood of } y} \right],$$

NPO

$$\mathcal{L}_{\text{NPO}}(\theta) = -\frac{2}{\beta} \mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_f} \left[\underbrace{\log \sigma \left(-\beta \log \frac{p(y|\mathbf{x}; \theta)}{p(y|\mathbf{x}; \theta)} \right)}_{\text{Decrease likelihood of } y} \right].$$

Two issues

- Unbounded forget loss & unclear optimization target

Unstable training &
gibberish output

$$\mathcal{L}_{\text{GA}}(\boldsymbol{\theta}) = -\mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_f} [-\log(p(y|\mathbf{X} = \mathbf{x}; \boldsymbol{\theta}))] = \mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_f} [\log(p(y|\mathbf{X} = \mathbf{x}; \boldsymbol{\theta}))].$$

Model collapse

$$\mathcal{L}_{\text{DPO}}(\boldsymbol{\theta}) = -\frac{1}{\beta} \mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_f, y^{idk} \sim \mathcal{D}_{idk}} = \left[\log \sigma \left(\underbrace{\beta \log \frac{p(y^{idk}|\mathbf{x}; \boldsymbol{\theta})}{p(y^{idk}|\mathbf{x}; \boldsymbol{\theta})}}_{\text{Increase likelihood of } y^{idk}} - \underbrace{\beta \log \frac{p(y|\mathbf{x}; \boldsymbol{\theta})}{p(y|\mathbf{x}; \boldsymbol{\theta})}}_{\text{Decrease likelihood of } y} \right) \right],$$

Two issues

- Unbounded forget loss & unclear optimization target

Unstable training &
gibberish output

$$\mathcal{L}_{\text{GA}}(\boldsymbol{\theta}) = -\mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_f} [-\log(p(y|\mathbf{X} = \mathbf{x}; \boldsymbol{\theta}))] = \mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_f} [\log(p(y|\mathbf{X} = \mathbf{x}; \boldsymbol{\theta}))].$$

Model collapse

$$\mathcal{L}_{\text{DPO}}(\boldsymbol{\theta}) = -\frac{1}{\beta} \mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_f, y^{\text{idk}} \sim \mathcal{D}_{\text{idk}}} \left[\log \sigma \left(\underbrace{\beta \log \frac{p(y^{\text{idk}}|\mathbf{x}; \boldsymbol{\theta})}{p(y^{\text{idk}}|\mathbf{x}; \boldsymbol{\theta})}}_{\text{Increase likelihood of } y^{\text{idk}}} - \underbrace{\beta \log \frac{p(y|\mathbf{x}; \boldsymbol{\theta})}{p(y|\mathbf{x}; \boldsymbol{\theta})}}_{\text{Decrease likelihood of } y} \right) \right],$$

- Under-representative retain loss

Only retain
knowledge in \mathcal{D}_r

$$\mathcal{L}_{\text{GD}}(\boldsymbol{\theta}) = \mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_r} [-\log(p(y|\mathbf{x}; \boldsymbol{\theta}))],$$

$$\mathcal{L}_{\text{KL}}(\boldsymbol{\theta}) = \mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}_w} \left[D_{\text{KL}} \left(p(y|\mathbf{x}; \boldsymbol{\theta}) \parallel p(y|\mathbf{x}; \boldsymbol{\theta}^{(0)}) \right) \right],$$

Two issues

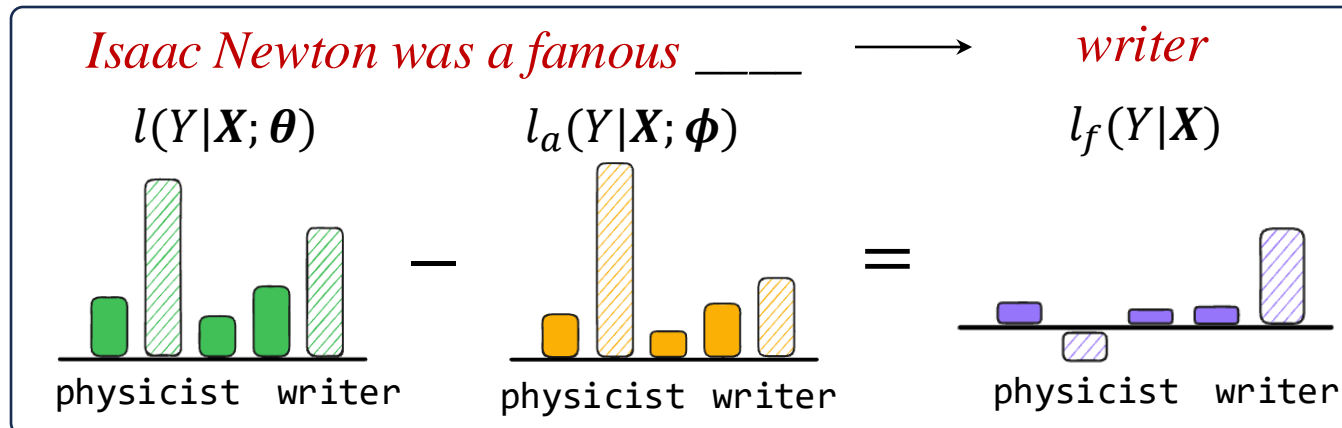
	<i>Query for forget documents</i>	<i>Query for retain documents</i>	<i>Query for knowledge not included in retain documents</i>
	When and where was Sir Isaac Newton born?	When and where was Aristotle born?	When and where was Geoffery Hinton born?
	Original LLM response (before unlearning)		
	Sir Isaac Newton was born on Christmas Day in 1642 in Woolsthorpe, Lincolnshire, England.	Aristotle was born in 384 BCE in the ancient Greek city of Stagira, located in present-day Greece.	Geoffrey Hinton was born on December 6, 1947. He was born in Wimbledon, England.
	Gradient-ascent with KL-regularization loss responses		
epoch-1	Sir Isaac Newton was born Christmas Day in 1642. He was born in Woolsthorpe, Lincolnshire, England.	Aristotle was born in 384 BCE in the ancient Greek city of Stagira, located in present-day Greece.	Geoffrey Hinton was born December 6, 1947. He was born in Wimbledon, England.
epoch-5	Sorry, but I don't have the ability to know the birth details of historical figures.	Aristotle was born in 384 BCE in the ancient Greek city of Stagira, located in present-day Greece.	Sorry, I don't know when or where Geoffrey Hinton was born.
epoch-10	Sorry Christmas Christmas Christmas Christmas Christmas Christmas Christmas . . .	Aristotle was born in 384 BCE Christmas Christmas Christmas Christmas . . .	I apologize Christmas Christmas Christmas Christmas Christmas Christmas . . .

Our solution: Unlearn from Logit Difference

- We seek an assistant LLM that remembers D_f , but no knowledge about D_r
- Then ensemble with original model to simulate unlearn in decoding

$$l_f(Y|\mathbf{X}) = l(Y|\mathbf{X}; \theta) - \alpha \cdot l_a(Y|\mathbf{X}; \phi),$$

ULD Inference on Query to Forget Knowledge

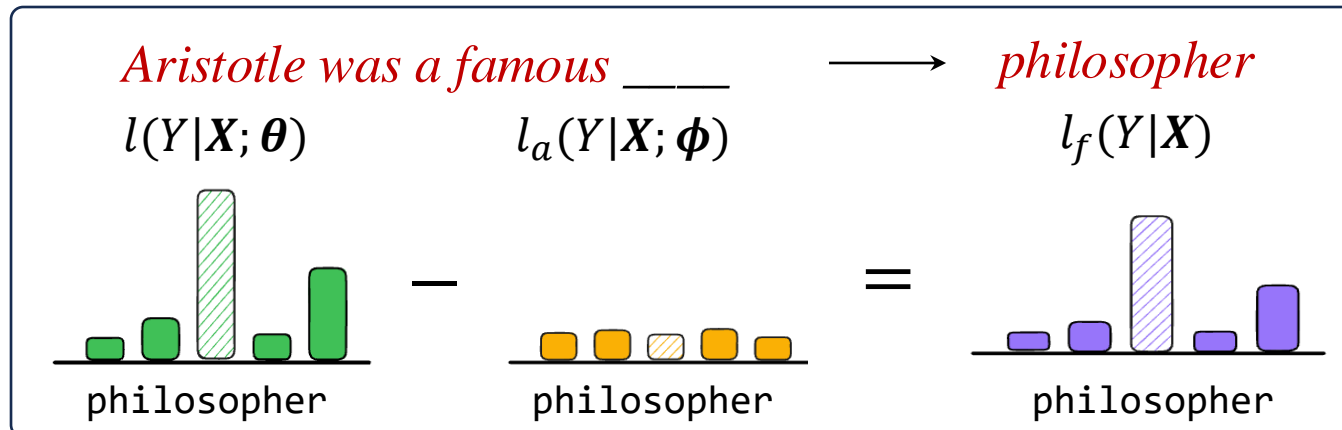


Our solution: Unlearn from Logit Difference

- We seek an assistant LLM that remembers D_f , but no knowledge about D_r
- Then ensemble with original model to simulate unlearn in decoding

$$l_f(Y|\mathbf{X}) = l(Y|\mathbf{X}; \boldsymbol{\theta}) - \alpha \cdot l_a(Y|\mathbf{X}; \boldsymbol{\phi}),$$

ULD Inference on Query to Retain Knowledge



Our solution: Unlearn from Logit Difference

- Training assistant with a well-defined objective:

$$\min_{\phi} \mathcal{L}(\phi) = \min_{\phi} \mathcal{L}_f(\phi) - \beta \mathcal{L}_r(\phi).$$

Remember \mathcal{D}_f

$$\mathcal{L}_f(\phi) = \mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}'_f} [\text{CE}(\text{softmax}(l_a(Y|\mathbf{X} = \mathbf{x}; \phi)); \delta(Y = y))],$$

Remain ignorant of \mathcal{D}_r

$$\mathcal{L}_r(\phi) = -\mathbb{E}_{\mathbf{x} \sim \mathcal{D}'_r} [\text{CE}(\text{softmax}(l_a(Y|\mathbf{X} = \mathbf{x}; \phi)); U(Y))],$$

Our solution: Unlearn from Logit Difference

- Training assistant with a well-defined objective:

$$\min_{\phi} \mathcal{L}(\phi) = \min_{\phi} \mathcal{L}_f(\phi) - \beta \mathcal{L}_r(\phi).$$

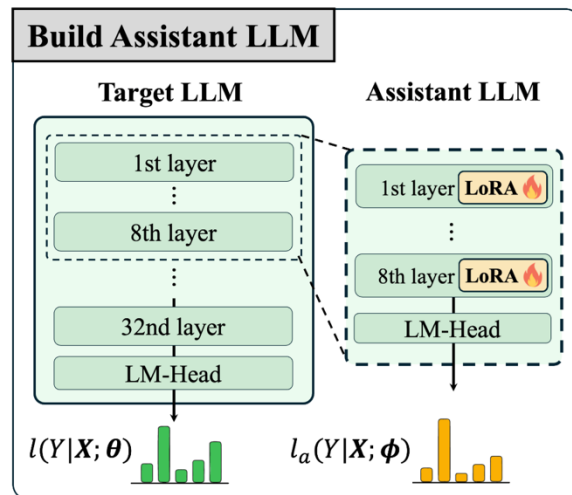
Remember \mathcal{D}_f

$$\mathcal{L}_f(\phi) = \mathbb{E}_{[\mathbf{x}, y] \sim \mathcal{D}'_f} [\text{CE}(\text{softmax}(l_a(Y|\mathbf{X} = \mathbf{x}; \phi)); \delta(Y = y))],$$

Remain ignorant of \mathcal{D}_r

$$\mathcal{L}_r(\phi) = -\mathbb{E}_{\mathbf{x} \sim \mathcal{D}'_r} [\text{CE}(\text{softmax}(l_a(Y|\mathbf{X} = \mathbf{x}; \phi)); U(Y))],$$

- An efficient training scheme



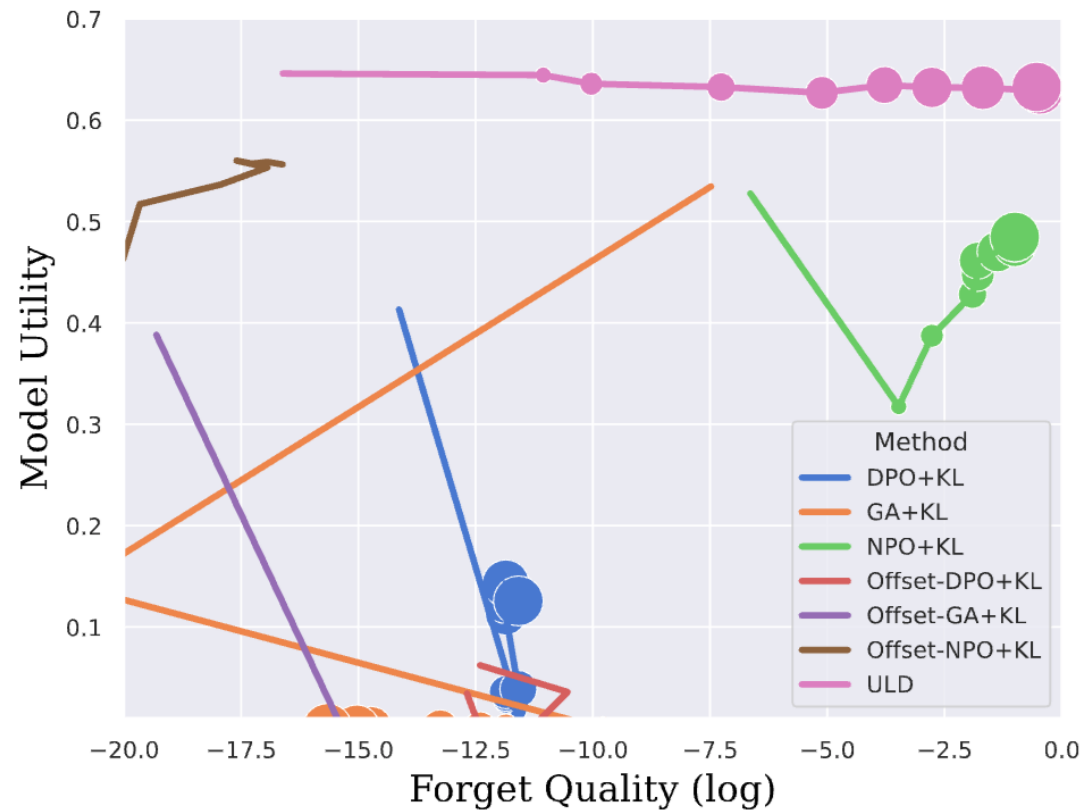
20M trainable parameter, 0.28% full model

Main result

TOFU: unlearn fictional author information

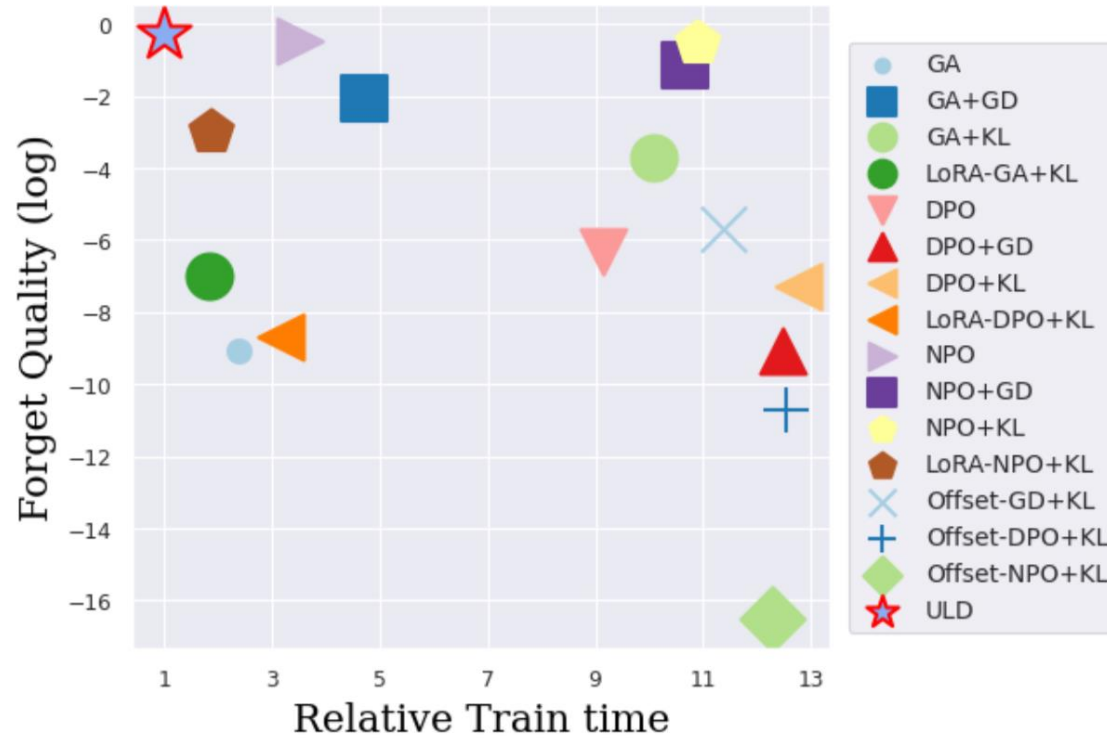
Method	TOFU-1%				TOFU-5%				TOFU-10%			
	Forget Perf.		Retain Perf.		Forget Perf.		Retain Perf.		Forget Perf.		Retain Perf.	
	<i>F.Q.</i> ↑	<i>R-L</i>	<i>M.U.</i> ↑	<i>R-L</i> ↑	<i>F.Q.</i> ↑	<i>R-L</i>	<i>M.U.</i> ↑	<i>R-L</i> ↑	<i>F.Q.</i> ↑	<i>R-L</i>	<i>M.U.</i> ↑	<i>R-L</i> ↑
Target LLM	1e-3	95.2	0.62	98.2	3e-16	97.3	0.62	98.2	2e-19	98.6	0.62	98.2
Retain LLM	1.0	37.6	0.62	98.5	1.0	39.3	0.62	98.1	1.0	39.8	0.62	98.2
GA	0.40	34.4	0.52	59.6	0.05	24.4	0.37	31.3	8e-10	0	0	0
GA+GD	0.27	30.5	0.53	58.9	0.11	19.5	0.33	28.9	9e-3	19.6	0.17	23.9
GA+KL	0.40	35.2	0.53	59.9	0.14	20.3	0.35	29.2	2e-4	12.1	0.05	18.6
DPO	0.27	4.09	0.58	55.2	1e-4	1.1	0.02	0.89	5e-7	0.7	0	0.72
DPO+GD	0.25	4.08	0.58	56.5	1e-7	1.2	0.02	0.84	8e-10	0.8	0	0.89
DPO+KL	0.26	4.18	0.58	55.6	4e-5	1.1	0.03	0.93	5e-8	0.7	0.03	0.81
NPO	0.66*	39.2	0.52	62.8	0.68	15.9	0.19	24.6	0.09	15.2	0.26	15.3
NPO+GD	0.58*	34.5	0.57	63.1	0.46	24.7	0.44	36.5	0.29	25.7	0.53	41.1
NPO+KL	0.52*	33.7	0.54	58.7	0.44	24.2	0.48	40.2	0.07	18.1	0.32	22.9
Offset-GA+KL	0.27	44.7	0.52	45.8	1e-4	1.2	0	0	2e-6	3.1	0.04	2.9
Offset-DPO+KL	0.13	3.8	0.12	19.1	2e-8	0	0	0	3e-9	1.3	0.02	1.4
Offset-NPO+KL	0.41	31.4	0.43	34.5	5e-10	37.3	0.59	40.9	4e-5	34.2	0.48	34.8
ULD	0.99	40.7	0.62	98.3	0.73	41.2	0.62	93.4	0.48	42.6	0.62	85.9

Training stability



Trajectory of Model utility versus forget quality (log) for different unlearning method on TOFU-10%

Training efficiency



*Log forget quality versus relative training time to ULD on TOFU-10%.
The top-left corner indicates better forget performance and efficiency.*

Data Usage Ablation

Method	Data config		Forget Perf.		Retain Perf.	
	\mathcal{D}'_f	\mathcal{D}'_r	$F.Q. \downarrow$	$R-L$	$M.U. \uparrow$	$R-L \uparrow$
Target LLM	-	-	2e-19	98.6	0.62	98.2
Retain LLM	-	-	1.0	39.8	0.62	98.2
GA+KL	✓	✓	4e-7	0	0	0
DPO+KL	✓	✓	7e-11	0	0	0
NPO+KL	✓	✓	1e-4	12.3	0.08	18.4
Offset-NPO+KL	✓	✓	6e-9	15.8	0.24	28.7
ULD	✗	✗	1e-7	13.7	0.53	34.1
ULD	✗	✓	1e-9	43.8	0.63	84.1
ULD	✓	✗	0.51	12.7	0.55	72.3
ULD	✓	✓	0.52	42.4	0.62	86.4

TOFU-10% unlearning performance for different methods using augmented Forget/Retain Data. The data augmentation is mostly useful for our method.

Thanks!



Paper



Code

Contact: jiabaoji@ucsb.edu