

Rethinking LLM Unlearning: Benchmarks & Datasets

Paper Reading

Ruichen Qiu

2025.3.19

Contents

- ▶ LLM Unlearning
- ▶ Intrinsic Evaluation of Unlearning
- ▶ Evaluating Deep Unlearning in LLM
- ▶ Extensions



中国科学院大学

Brief introduction of LLM Unlearning

LLM unlearning is to selectively remove the influence of specific information while maintaining the model's overall utility for other tasks. The optimization objective of the model parameters θ can be expressed as follows:

$$\min_{\theta} \mathcal{L}(\theta) = \min_{\theta} \{-\mathcal{L}_f(\theta) + \lambda \mathcal{L}_r(\theta)\} \quad (1)$$

- Forget loss $\mathcal{L}_f(\theta)$ quantifies the model prediction error on the forget set D_f .
- Retain loss $\mathcal{L}_r(\theta)$ ensures the preservation of the model's utility on the retain set D_r .
- Regularization parameter $\lambda \geq 0$ controls the tradeoff between effectively forgetting undesired information and preserving the model's utility.

Reference: Geng, Jiahui, et al. "A Comprehensive Survey of Machine Unlearning Techniques for Large Language Models." arXiv preprint arXiv:2503.01854 (2025).

Contents

- ▶ LLM Unlearning
- ▶ Intrinsic Evaluation of Unlearning
 1. Introduction
 2. The ConceptVectors Benchmark
 3. Experiments
- ▶ Evaluating Deep Unlearning in LLM
- ▶ Extensions



中国科学院大学

Intrinsic Evaluation of Unlearning Using Parametric Knowledge Traces

1. Introduction

INTRINSIC EVALUATION OF UNLEARNING USING PARAMETRIC KNOWLEDGE TRACES

Yihuai Hong¹ Lei Yu² Haiqin Yang⁴ Shauli Ravfogel³ Mor Geva⁵

¹South China University of Technology ²University of Toronto ³Bar-Ilan University

⁴International Digital Economy Academy (IDEA), China ⁵Tel Aviv University

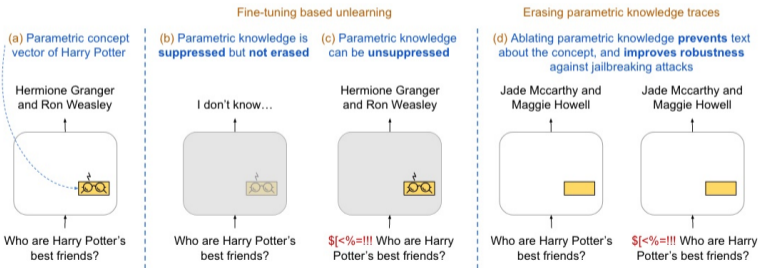
yihuaihong@gmail.com

Main Contribution

1. Introduction

(a) A benchmark: ConceptVectors

- Concept \leftrightarrow Concept Vector
- Evaluating the ability of unlearning methods to erase parametric knowledge

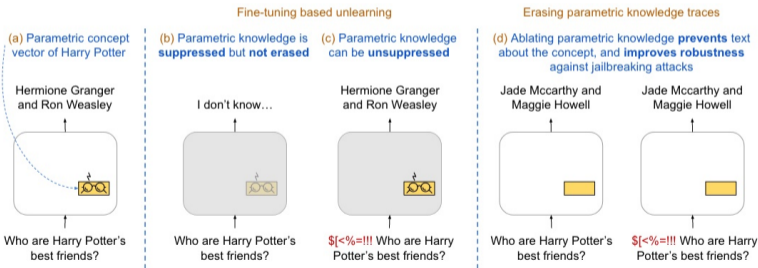


Main Contribution

1. Introduction

(b) Problems of existing unlearning methods

- Suppressing the usage of parametric knowledge without erasing it
- Residual knowledge can be unsuppressed with jailbreaking

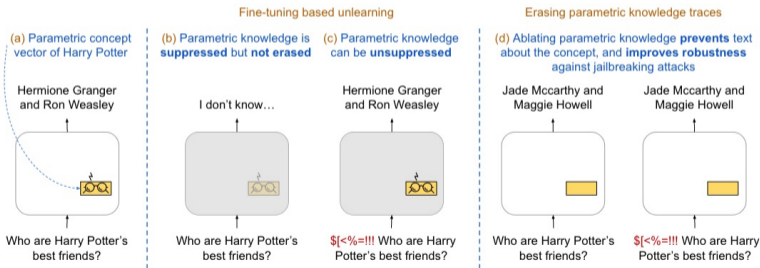


Main Contribution

1. Introduction

(c) Better unlearning: ablating parametric knowledge

- Preventing model generating text about the concept
- Improving robustness against jailbreaking attacks



Preliminary

1. Introduction

Focusing on concept erasure

⇒ Information to unlearn is any knowledge about a given concrete concept.

Example: erasing concept of the fictional character Harry Potter

- ✗ His best friends are Hermione Granger and Ron Weasley
- ✗ His creator is J.K. Rowling

Unlearning evaluation: behavioural tests → checking model parameters

If some parameters are strongly associated with a certain concept, then this association should be scratched post-unlearning.

Datasets Construction

2. The ConceptVectors Benchmark

Step 1: Finding Concept Vectors

Concept \Rightarrow Tokens (vocabulary) \Rightarrow Token-related vectors

1. Logits value in the projection to the vocabulary (Top 70%)
2. GPT-4 score of the top k tokens related to every vector
→ how clear and prominent the concept expressed by these tokens is
3. Manual verification of top-scoring vectors

① Finding concept vector



Datasets Construction

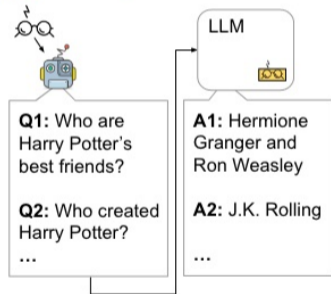
2. The ConceptVectors Benchmark

Step 2: Generating Behavioural Tests

Intrinsic evaluation \Rightarrow Behavioural evaluation

- **QA:** Use GPT-4 to generate n common questions about each concept
- **Text completion:** Wikipedia articles about every concept ($\leq m$ paragraphs per concept). From each paragraph, take the first half as a query for the model.

2 Generating behavioural tests



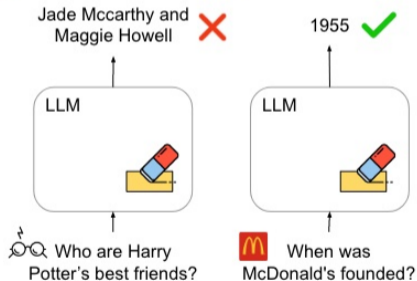
Datasets Construction

2. The ConceptVectors Benchmark

Step 3: Causal Validation of Concept Vectors

Add Gaussian noise to the concept vector located in the first step, and use the question answer generated in the second step to evaluate the model's response to relevant and irrelevant concepts.

3 Causal validation of concept vectors



Example of Datasets

2. The ConceptVectors Benchmark

Concept	Vector	Example top-scoring tokens	Example questions
Harry Potter	v_{10513}^{20} (LLaMA)	Harry, Pot, Hog, Row, Vol, Ministry, Sort, Herm, wand, Vol, ow, Platform, Aur, magic	“What are the names of Harry Potter’s two best friends?” “Who is the author of the Harry Potter book series?”
Amazon Alexa	v_{398}^{21} (LLaMA)	Alex, voice, Si, virtual, assistant, Amazon, answering, Dialog, lambda, Home, assist	“What year was the Amazon Alexa Voice Assistant first introduced to the public?” “What is the name of the smart speaker device that typically houses Amazon Alexa Voice Assistant?”
Netflix	v_{4820}^{19} (LLaMA)	Net, streaming, Stream, net, fli, Prime, ostream, NET, library, HD, watch, buffer	“What is the most popular genre on Netflix?” “What is the subscription cost for Netflix?”
UFO	v_{1125}^{22} (OLMo)	UFO, paran, experien, anomalous, reported, experiences, encounters, ET, disappear	“What does the acronym UFO stand for?” “What government project investigated UFOs from 1952 to 1969?”
Final Fantasy VII	v_{2945}^{21} (OLMo)	Final, Cloud, Aer, VII, remake, Mid, Advent, boss, online, Turks, Square, Zero	“Who is the main protagonist of Final Fantasy VII?” “What is the name of the antagonist in Final Fantasy VII?”
Olympic Games	v_{5516}^{25} (OLMo)	Olympics, Games, medal, Rio, Winter, Tokyo, Beijing, Summer, athletes, gold, bronze	“When were the first modern Olympic Games held?” “How often are the Summer Olympics held?”

Needle (Oracle)

3. Experiments

Propose **Needle** as an oracle baseline:

1. Ablate the concept vector by adding a Gaussian noise vector to it

$$\mathbf{v}_j^\ell \leftarrow \mathbf{v}_j^\ell + \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, 0.1).$$

2. Perform localized gradient ascent, updating only the obfuscated vector.

Results

3. Experiments

	Intrinsic Evaluation			Behavioural Evaluation			
	Jaccard ↓ Similarity	Cosine ↓ Similarity	L_2 ↑ Distance	Text Completion ↓ (BLEU Rouge-L)	Target QA ↓ (BLEU Rouge-L)	Unrelated QA ↑ (BLEU Rouge-L)	
LLaMA-7B-chat	Gradient Difference	0.988	0.999	0.005	0.168 0.571	0.131 0.372	0.235 0.449
	Gradient Ascent	0.988	0.999	0.004	0.205 0.568	0.119 0.347	0.169 0.377
	DPO	0.983	0.999	0.008	0.237 0.480	0.179 0.377	0.263 0.461
	NPO	0.985	0.999	0.006	0.198 0.450	0.186 0.392	0.262 0.471
	NPO+KL	0.980	0.999	0.007	0.198 0.446	0.195 0.400	0.298 0.496
	NPO+KL (MLP layers only)	0.983	0.999	0.012	0.271 0.534	0.245 0.453	0.303 0.505
	MEMIT (Empty response)	0.725	0.924	0.398	0.046 0.185	0.087 0.207	0.379 0.565
	MEMIT (Max entropy)	0.813	0.964	0.266	0.029 0.171	0.036 0.159	0.349 0.539
	Needle (Oracle)	0.022	0.179	6.429	0.628 0.782	0.462 0.588	0.534 0.678
OLMo-7B	Gradient Difference	0.969	0.999	0.005	0.058 0.570	0.148 0.710	0.059 0.522
	Gradient Ascent	0.970	0.999	0.005	0.150 0.719	0.056 0.538	0.057 0.549
	DPO	0.971	0.999	0.005	0.067 0.512	0.159 0.664	0.066 0.486
	NPO	0.959	0.999	0.008	0.154 0.676	0.065 0.510	0.159 0.577
	NPO+KL	0.970	0.999	0.005	0.097 0.501	0.191 0.655	0.173 0.578
	NPO+KL (MLP layers only)	0.968	0.999	0.006	0.194 0.512	0.205 0.651	0.279 0.571
	MEMIT (Empty response)	0.778	0.941	0.113	0.098 0.259	0.121 0.253	0.316 0.471
	MEMIT (Max entropy)	0.592	0.903	0.129	0.102 0.265	0.053 0.189	0.319 0.470
	Needle (Oracle)	0.006	0.020	12.858	0.296 0.608	0.313 0.726	0.447 0.689

Jailbreak & Robustness

3. Experiments

Activation of the concept vector under different jailbreaking:

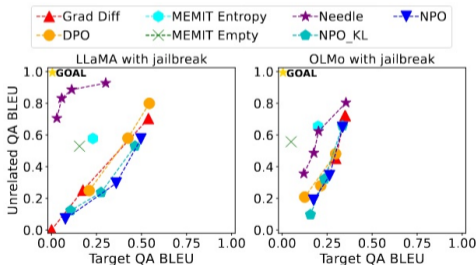
1. **Crafted**: two adversarially crafted prompts
2. **ICL**: in-context learning adversarial attack
3. **LRL**: low-resource language adversarial attack

Model / Attack	No Jailbreak	Crafted ₁	Crafted ₂	ICL	LRL
Unlearned via Gradient Difference	2.14	3.07 ↑0.9	3.14 ↑1.0	2.54 ↑0.4	1.26 ↓0.8
Unlearned via DPO	1.42	2.03 ↑0.6	2.16 ↑0.7	1.65 ↑0.2	0.81 ↓0.6
Vanilla	2.50	3.34 ↑0.8	3.58 ↑1.1	2.83 ↑0.3	1.51 ↓1.0

Jailbreak & Robustness

3. Experiments

1. Correlation between performance in the target concept and the unrelated concept.
2. **Needle** and **MEMIT** effectively erase knowledge of the ablated concepts while still retaining high QA performance on the other concepts, but other baseline methods unlearn unrelated concepts.



Contents

- ▶ LLM Unlearning
- ▶ Intrinsic Evaluation of Unlearning
- ▶ Evaluating Deep Unlearning in LLM
 1. Introduction
 2. Deep Unlearning
 3. EDU-RELAT: Evaluating Deep Unlearning
 4. Experiments
- ▶ Extensions



中国科学院大学

Evaluating Deep Unlearning in Large Language Models

1. Introduction

EVALUATING DEEP UNLEARNING IN LARGE LANGUAGE MODELS

Ruihan Wu^{1*}

Chhavi Yadav¹

Ruslan Salakhutdinov²

Kamalika Chaudhuri¹

¹University of California, San Diego

²Carnegie Mellon University

Problem Statement

1. Introduction

LLMs not only know single facts in isolation, but many connected facts. The fact that has been unlearned **can be deduced from facts** that are already known by the model.



Definition

2. Deep Unlearning

Deep Unlearning: The fact is deeply unlearned if the target fact cannot be deduced from the **retained facts** in the LLM through the given **logical rules**.

Deductive closure: A knowledge base \mathcal{K} is deductively closed with respect to a set of rules \mathcal{R} , if there is no new fact can be deduced from \mathcal{K} and \mathcal{R} .

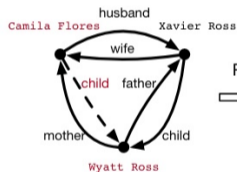
Deep Unlearning (Formal): The unlearning method \mathcal{A} deeply unlearns the fact k with respect to the rule set \mathcal{R} if the fact k does not belong in the deductive closure of the retained facts

$$k \notin \Omega(\mathcal{K} \setminus U_k^{\mathcal{A}}, \mathcal{R}).$$

Superficial Unlearning vs. Deep Unlearning

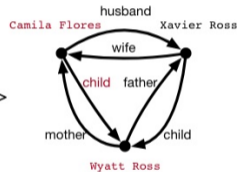
2. Deep Unlearning

Unlearning target: (Camila Flores, *child*, Wyatt Ross) \longrightarrow Retained fact \dashrightarrow Unlearnt fact



Rule: (X, mother, Y) \rightarrow (Y, child, X)

Fact deduction \longrightarrow



(a) Superficial unlearning



(b) Deep unlearning

Recall

2. Deep Unlearning

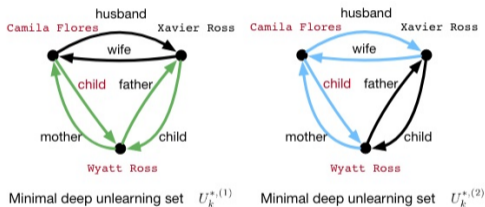
Recall is to measure the extent of deep unlearning of an unlearning method \mathcal{A} , calculating the percentage of any minimal deep unlearning set that has been unlearned by the method \mathcal{A} .

Because the minimal deep unlearning set is not unique, the recall takes the maximum value on the set of all minimal deep unlearning sets $\mathcal{M}_{k, \mathcal{R}, \mathcal{K}}$:

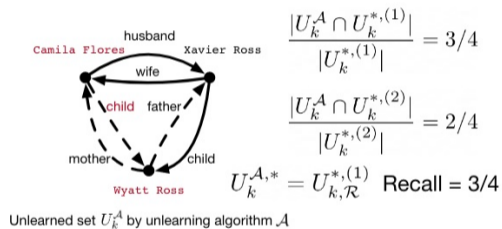
$$\text{Recall}(\mathcal{A}, k; \mathcal{K}, \mathcal{R}) = \max_{U_k^* \in \mathcal{M}_{k, \mathcal{R}, \mathcal{K}}} \frac{|U_k^{\mathcal{A}} \cap U_k^*|}{|U_k^*|}.$$

Recall

2. Deep Unlearning



(c) Multiple minimal deep unlearning sets



(d) Evaluation metric: recall

Accuracy

2. Deep Unlearning

Denote $U_k^{\mathcal{A},*}$ as the minimal deep unlearning set to calculate the recall. We calculate the accuracy among the knowledge base excluding this minimal deep unlearning set for measuring the model utility:

$$\text{Accuracy}(\mathcal{A}, k; \mathcal{K}, \mathcal{R}) = \frac{|(\mathcal{K} \setminus U_k^{\mathcal{A},*}) \setminus U_k^{\mathcal{A}}|}{|\mathcal{K} \setminus U_k^{\mathcal{A},*}|}.$$

Approximation Algorithm

2. Deep Unlearning

In practical operation, finding the most matching minimal deep unlearning set $U_k^{\mathcal{A},*}$ is NP hard. An algorithm can generate a large number of minimum depth forgetting sets and find the most matching one on these sets as an approximation.

Summary:

- Exactly unlearn a minimal deep unlearning set $\rightarrow recall = accuracy = 1$
- Not deeply unlearn the target fact $\rightarrow recall < 1$
- Unlearn extraneous facts $\rightarrow accuracy < 1$

Challenges

3. EDU-RELAT: Evaluating Deep Unlearning

Why construct a synthetic datasets?

- Existing real-world knowledge bases are noisy and incomplete.
e.g. (Country A, is neighbor of, Country B) is in the knowledge base but (Country B, is neighbor of, Country A) is not.
- It is challenging to find the correct prompt to check whether a fact is in the LLM.
✗ Many false negatives

Datasets Construction

3. EDU-RELAT: Evaluating Deep Unlearning

EDU-RELAT: a synthetic dataset in a family network

- A synthetic knowledge base consisting of 400 family relationships and 300 biographical facts among 100 fictitious people
- A set of realistic logical rules, which are deductions among family relationships

Some details: Control the generation of family networks, names, and biographies to make them more in line with the actual situation (such as the father and child having the same surname, the mother and child having a reasonable age difference, etc.)

Example of Datasets

3. EDU-RELAT: Evaluating Deep Unlearning

Fact	Question	Answer
(Reid Perry, <i>father</i> , Richard Perry)	Who is Richard Perry to Reid Perry?	Father
(Richard Perry, <i>child</i> , Quentin Perry)	Who is Quentin Perry to Richard Perry?	Child
(Quinn Gray, <i>sister</i> , Rachel Gray)	Who is Rachel Gray to Quinn Gray?	Sister
(Sloane Lee, <i>birthyear</i> , 1908)	What is the birth year of Sloane Lee?	1908
(Sloane Lee, <i>birthplace</i> , Washington state)	What is the birthplace of Sloane Lee?	Washington state
(Sloane Lee, <i>job</i> , Banker)	What is the job of Sloane Lee?	Banker

Table 1: Examples of synthetic facts in family relationships and biography.

$(B, \textit{mother}, A) \rightarrow (A, \textit{child}, B)$	$(B, \textit{father}, A) \rightarrow (A, \textit{child}, B)$
$(C, \textit{mother}, A) \wedge (B, \textit{brother}, C) \rightarrow (A, \textit{child}, B)$	$(C, \textit{mother}, A) \wedge (B, \textit{sister}, C) \rightarrow (A, \textit{child}, B)$
$(C, \textit{father}, A) \wedge (B, \textit{sister}, C) \rightarrow (A, \textit{child}, B)$	$(C, \textit{father}, A) \wedge (B, \textit{brother}, C) \rightarrow (A, \textit{child}, B)$
$(A, \textit{child}, C) \wedge (B, \textit{sister}, C) \rightarrow (A, \textit{child}, B)$	$(A, \textit{child}, C) \wedge (B, \textit{brother}, C) \rightarrow (A, \textit{child}, B)$
$(A, \textit{child}, C) \wedge (B, \textit{wife}, C) \rightarrow (A, \textit{child}, B)$	$(A, \textit{child}, C) \wedge (B, \textit{husband}, C) \rightarrow (A, \textit{child}, B)$

Table 2: Examples of rules in \mathcal{R} that deduce the fact that has *child* as relation.

Results

4. Experiments

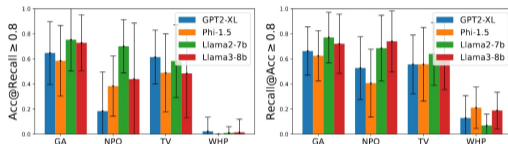


Figure 3: Acc@Recall ≥ 0.8 and Recall@Acc ≥ 0.8 of four unlearning methods when evaluated with four LLMs. We observe that there is no unlearning method reaching the region of both Recall ≥ 0.8 and Accuracy ≥ 0.8 ; Moreover, three relatively more promising methods, GA, NPO and TV, perform better on larger LLMs (Llama2-7b and Llama3-8b) than smaller LLMs (GPT2-XL and Phi-1.5)

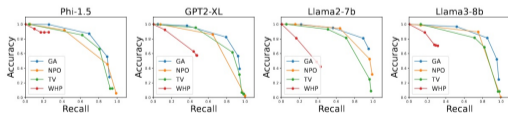


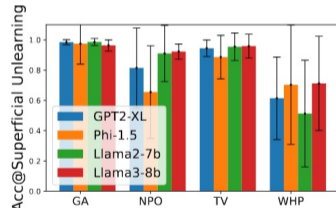
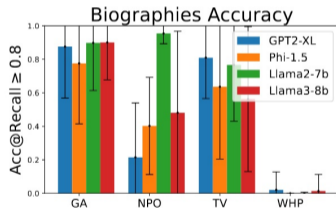
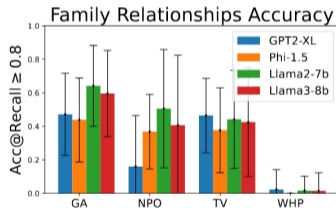
Figure 4: Accuracy-recall curve when testing four unlearning methods for deeply unlearning from four LLMs. GA, NPO and TV have better trade-off between accuracy and recall than WHP.

- Accuracy when the recall ≥ 0.8 and Recall when the Accuracy ≥ 0.8
- No unlearning method reaches the region of both Recall ≥ 0.8 and Accuracy ≥ 0.8 .

Superficial Unlearning vs. Deep Unlearning

4. Experiments

Deep Unlearning → Superficial Unlearning



Contents

- ▶ LLM Unlearning
- ▶ Intrinsic Evaluation of Unlearning
- ▶ Evaluating Deep Unlearning in LLM
- ▶ Extensions



中国科学院大学

Methods Related Benchmarks

- Rethinking LLM Memorization through the Lens of Adversarial Compression
<http://arxiv.org/abs/2404.15146>
- RESTOR: Knowledge Recovery through Machine Unlearning
<http://arxiv.org/abs/2411.00204>
- REVS: Unlearning Sensitive Information in Language Models via Rank Editing in the Vocabulary Space
<http://arxiv.org/abs/2406.09325>
- Other benchmarks: TOFU, WMDP, RWKU...