

LLM Unlearning on Noisy Forget Sets: A Study of Incomplete, Rewritten, and Watermarked Data





Changsheng Wang¹, Yihua Zhang¹, Jinghan Jia¹, Dennis Wei², Pin-Yu Chen², Sijia Liu^{1,2} ¹Michigan State University, ²IBM Research



Forget Data "Noise" in LLM Unlearning



Regulatory peptides control various important physiological processes ranging from fertilisation.



Regulatory peptides **** various *** physiological *** *** ranging *** fertilisation.



Regulatory peptides play key roles in a wide range of physiological processes, including fertilization.



Regulatory peptides are involved in diverse physiological functions, from fertilization and beyond.

Figure 1. Different potential perturbation types applied to the original data.

Motivation: "Noisy" (non-adversarial) forget data present significant challenges to the robustness of unlearning.

> Incomplete Forget Data vs. Unlearning

Only partial information is available for unlearning. We define $Mask_{\delta}(\mathbf{x})$ as a function that randomly masks δ (%) of tokens in each forget sample $\mathbf{x} \in \mathcal{D}_f$, producing a noisy forget set with uniformly sampled masked positions.

$$\mathcal{D}_{f}' = \{ Mask_{\delta}(\mathbf{x}) \mid \forall \mathbf{x} \in \mathcal{D}_{f} \}$$

Tolerance of unlearning to masking ratio (<= 30%)

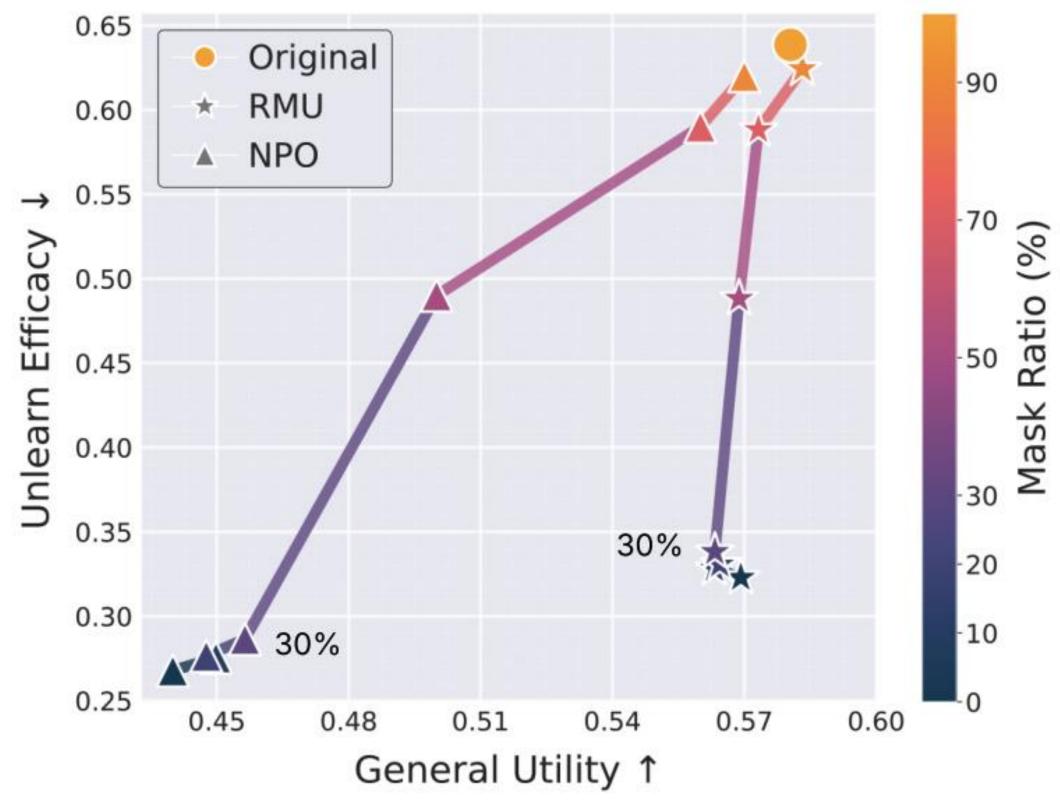


Figure 2. Impact of masking ratio on unlearning performance across two representative unlearning methods, NPO and RMU, applied to the Zephyr-7b-beta model on the WMDP dataset. Unlearn Efficacy is WMDP accuracy.

Rewritten Forget Data vs. Unlearning

Rewrite(·) denote a rewriting function that generates a paraphrased variant of a forget sample x while preserving its original semantics.

$$\mathcal{D}_{f}' = \{ \text{Rewrite}(\mathbf{x}) \mid \forall \mathbf{x} \in \mathcal{D}_{f} \}$$

Forget data type	UE ↓	UT↑
No unlearning	0.6386	0.5805
Clean data	0.3229	0.5692
Mask	0.3382	0.5632
Rewrite	0.3142	0.5680
WM (KGW)	0.3134	0.5694
WM (SynID)	0.3221	0.5684

Table 1. Performance of RMU - unlearning on perturbed forget data _ using Zephyr-7B-beta. Comparison of unlearning efficacy and general utility on the WMDP benchmark under various forget data conditions. UE means unlearn efficacy and UT is MMLU accuracy.

Watermarked Forget Data vs. Unlearning

Watermark(·) denotes the output of a watermarkenabled LLM decoding process for input x.

$$\mathcal{D}_{f}' = \{ \text{Watermark}(\mathbf{x}) \mid \forall \mathbf{x} \in \mathcal{D}_{f} \}$$

	These peptides often act as signaling	Watermark	RMU	NPO				
KGW (δ = 2)	molecules, coordinating responses to	Strength	UE ↓ UT ↑	UE ↓ UT ↑				
	internal and external stimuli. The study of	Original Model	0.6386 0.5805	0.6386 0.5805				
	regulatory peptides involves understanding their synthesis, processing, and function, as	Original Data	0.3229 0.5692	0.2603 0.4436				
	well as their interactions with receptors .	Logits-based Watermarking (KGW)						
KGW (δ = 6)	Peptids is like very important for signal	$\delta = 2$	0.3134 0.5694	0.2765 0.4521				
	thing. They making something inside and	$\delta = 4$	0.3652 0.5631	0.3124 0.4675				
	outside. Also they do many job like job	$\delta = 6$	0.3764 0.5461	0.3265 0.4613				
	of signal and another job. People study for Sampling-based Watermarking (SynthID)							
	why peptides happen and make and also how	m = 2	0.3201 0.5673	0.2675 0.4498				
	make and then doing thing with receptor or	m=4	0.3221 0.5684	0.2641 0.4501				
	other stuff. Function is also something about.	m = 6	0.3465 0.5512	0.2945 0.4598				

with representative stronger watermark signal.

Table 2. Watermark examples Table 3. Unlearning performance KGW under different watermarking watermarking method. For KGW, strengths. This table reports the tokens highlighted in red belong to unlearning performance of two the red list, and those in green representative unlearning methods, belong to the green list. Higher RMU and NPO, applied to the proportion of red tokens reflects a Zephyr-7b-beta model on the WMDP.

Other Experiment Results Highlights

Unlearning performance under perturbed forget data.

	UE			UT	
Forget data type	VerbMem (↓)	KnowMem (↓)	PrivLeak (→ 0)	KnowMem (†)	
No unlearning	99.80	59.40	-57.50	66.90	
Clean data	0.00	1.18	-42.07	57.19	
Mask	0.05	0.33	-49.36	55.31	
Rewrite	0.06	0.00	-53.43	50.73	
WM(KGW)	0.12	1.00	-53.51	56.92	
WM(SynthID)	0.05	1.13	-48.65	56.42	

Table 4. Unlearning performance of NPO on MUSE-Books using ICLM-7B under various forget data perturbations.

Analyzing error set overlap to assess unlearning robustness.

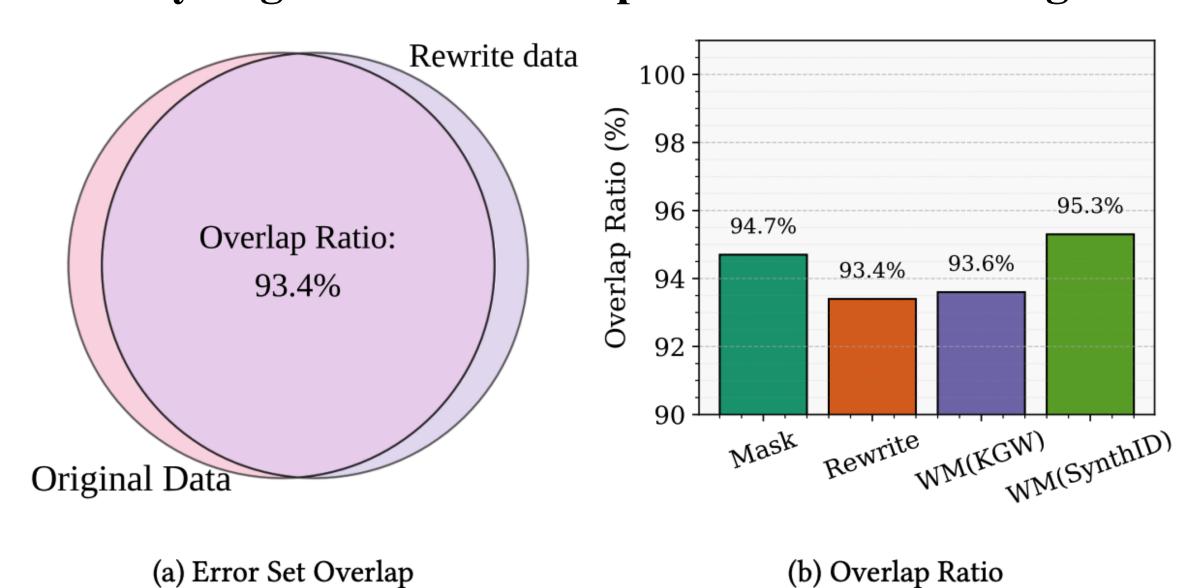


Figure 3. Performance consistency of unlearning error rates under perturbed forget data. (a) Venn diagram showing the overlap in incorrectly answered WMDP questions between models unlearned with original and rewritten forget data. (b) Overlap ratios between the error sets of models unlearned with various perturbed forget sets.

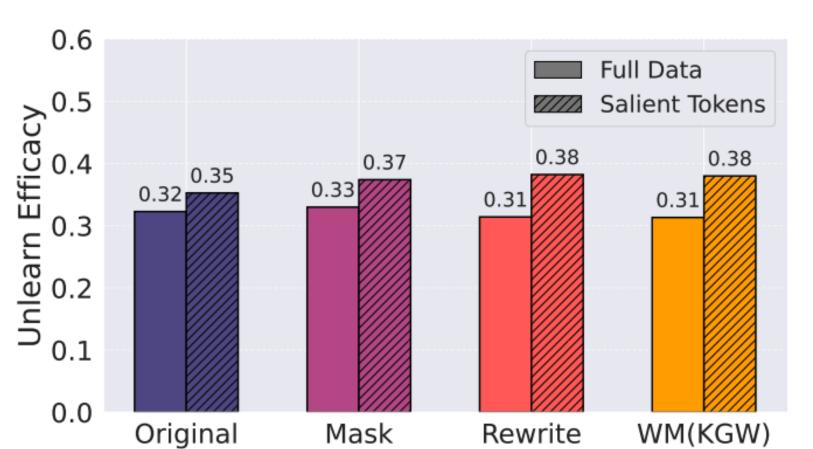


Figure Performance comparison RMU of unlearning on WMDP using Zephyr-7b-beta with full data salient tokens across original, mask, rewrite, and WM. Salient tokens achieve efficacy close to full data.