

# ***Learning to Unlearn while Retaining: Combating Gradient Conflicts in Machine Unlearning***

Gaurav Patel and Qiang Qiu



Elmore Family School of Electrical  
and Computer Engineering



# Context & Problem Statement

- **Machine Unlearning (MU):**

- Developed to comply with privacy rights like the “right to be forgotten” (GDPR).
- *Exact unlearning* = Retraining from scratch without the unwanted data—but that is very computationally expensive.
- *Approximate* unlearning methods aim to remove data influence more **efficiently**, by finetuning the initial model.

- **Competing Objectives:**

- In practice, two objectives:
  - **Forget** specific data (reduce the model’s reliance or memory of it).
  - **Retain** performance on remaining data.
- **Contradictory gradient** updates → can cause gradient conflict, slow or prevent proper unlearning.



# Preliminary

- **Standard Machine Unlearning Objective:**

- Typical approach combine the losses:

$$\min_{\theta} \left[ \underbrace{\mathcal{L}_r(\theta; \mathcal{D}_r)}_{\text{Retain Loss}} + \lambda \underbrace{\mathcal{L}_f(\theta; \mathcal{D}_f)}_{\text{Forget Loss}} \right]$$

- $\mathcal{L}_r$  (retain loss) ensures performance (utility) on the *retain set*  $\mathcal{D}_r$ .
    - $\mathcal{L}_f$  (forget loss) pushes the model to “unlearn” the forget set  $\mathcal{D}_f$ .
    - $\lambda$  balances the two objectives.

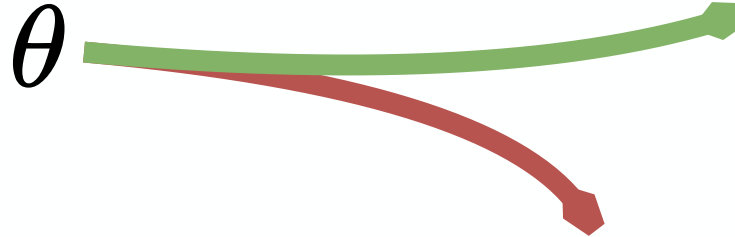
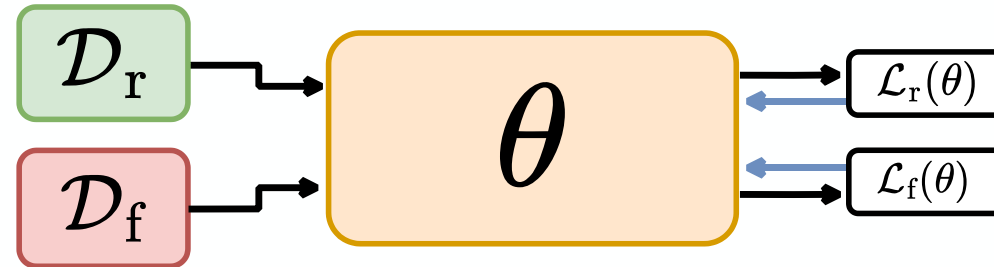
- **Issue:** When  $\nabla \mathcal{L}_r$  and  $\nabla \mathcal{L}_f$  point in opposing directions, updates can *cancel* each other’s effect.



# Our Proposition

## Naïve Unlearning

A.



# Our Proposition

## Learning to Unlearn while Retaining (LUR)

- **LUR Formulation:**

- **Key Idea:** Treat forgetting as the higher-level objective and retaining as the lower-level objective in a two-step update.

- **Implementation:**

- **Intermediate step:** Update  $\theta$  a little (with step  $\alpha$ ) using  $\mathcal{L}_r$ .

$$\theta' = \theta - \alpha \nabla \mathcal{L}_r(\theta)$$

- **Final Update:** Evaluate  $\mathcal{L}_f$  at  $\theta'$  (which accounts for the “retain” direction) and update  $\theta$ .

Formally:

$$\min_{\theta} [\mathcal{L}_r(\theta; \mathcal{D}_r) + \mathcal{L}_f(\theta'; \mathcal{D}_f)]$$



# Our Proposition

## Gradient Analysis

$$\mathcal{L}_r(\theta) + \mathcal{L}_f(\theta') = \mathcal{L}_r(\theta) + \mathcal{L}_f(\theta - \alpha \nabla \mathcal{L}_r(\theta))$$

- Gradient with respect to  $\theta$

$$\begin{aligned} \mathbf{g}_\theta &= \nabla \mathcal{L}_r(\theta) + \nabla \mathcal{L}_f(\theta') \cdot \frac{\partial \theta'}{\partial \theta} \\ &= \nabla \mathcal{L}_r(\theta) + \nabla \mathcal{L}_f(\theta') \cdot \frac{\partial (\theta - \alpha \nabla \mathcal{L}_r(\theta))}{\partial \theta} \\ &= \nabla \mathcal{L}_r(\theta) + \nabla \mathcal{L}_f(\theta') \cdot (I - \alpha \nabla^2 \mathcal{L}_r(\theta)) . \end{aligned}$$



# Our Proposition

## Gradient Analysis

- Gradient with respect to  $\theta$

$$\mathbf{g}_\theta = \nabla \mathcal{L}_r(\theta) + \overbrace{\left( \nabla \mathcal{L}_f(\theta) + \nabla^2 \mathcal{L}_f(\theta) \underbrace{(\theta' - \theta)}_{=-\alpha \nabla \mathcal{L}_r(\theta)} + \underbrace{\mathcal{O}(\|\theta' - \theta\|^2)}_{=\mathcal{O}(\alpha^2)} \right)}^{=\nabla \mathcal{L}_f(\theta')} (I - \alpha \nabla^2 \mathcal{L}_r(\theta))$$

$$\Rightarrow \mathbf{g}_\theta = \nabla \mathcal{L}_r(\theta) + \nabla \mathcal{L}_f(\theta) + \nabla^2 \mathcal{L}_f(\theta) (\theta' - \theta) - \alpha \nabla^2 \mathcal{L}_r(\theta) \nabla \mathcal{L}_f(\theta) + \mathcal{O}(\alpha^2)$$

$$(\text{using } \theta' - \theta = -\alpha \nabla \mathcal{L}_r(\theta))$$

$$\Rightarrow \mathbf{g}_\theta = \nabla \mathcal{L}_r(\theta) + \nabla \mathcal{L}_f(\theta) - \alpha \nabla^2 \mathcal{L}_f(\theta) \nabla \mathcal{L}_r(\theta) - \alpha \nabla^2 \mathcal{L}_r(\theta) \nabla \mathcal{L}_f(\theta) + \mathcal{O}(\alpha^2)$$

$$\Rightarrow \mathbf{g}_\theta = \nabla \mathcal{L}_r(\theta) + \nabla \mathcal{L}_f(\theta) - \underbrace{\alpha \nabla \left( \nabla \mathcal{L}_r(\theta) \cdot \nabla \mathcal{L}_f(\theta) \right)}_{\text{Gradient Product}} + \mathcal{O}(\alpha^2).$$

$$\theta \leftarrow \theta - \eta \mathbf{g}_\theta$$



# Our Proposition

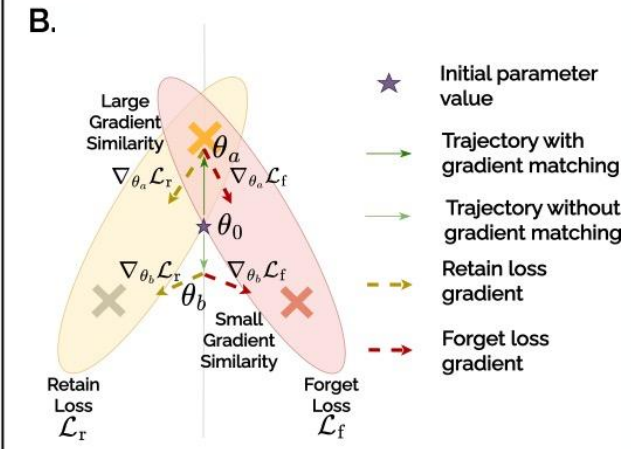
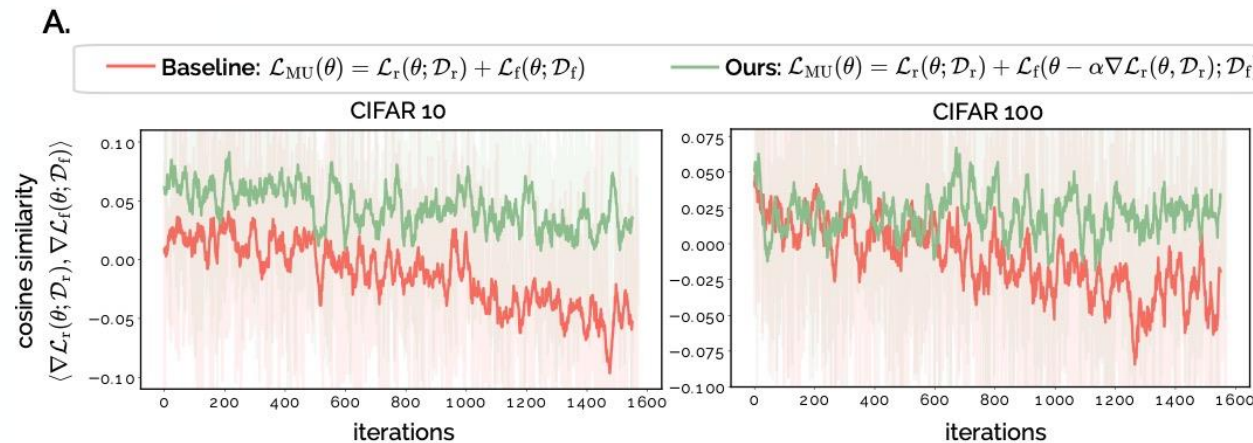
## Insight

### Why does it help?

- The model “anticipates” how forgetting will affect retention and *implicitly* promotes alignment of  $\nabla \mathcal{L}_r$  with  $\nabla \mathcal{L}_f$
- Mathematically shown via **implicit gradient product regularization**:

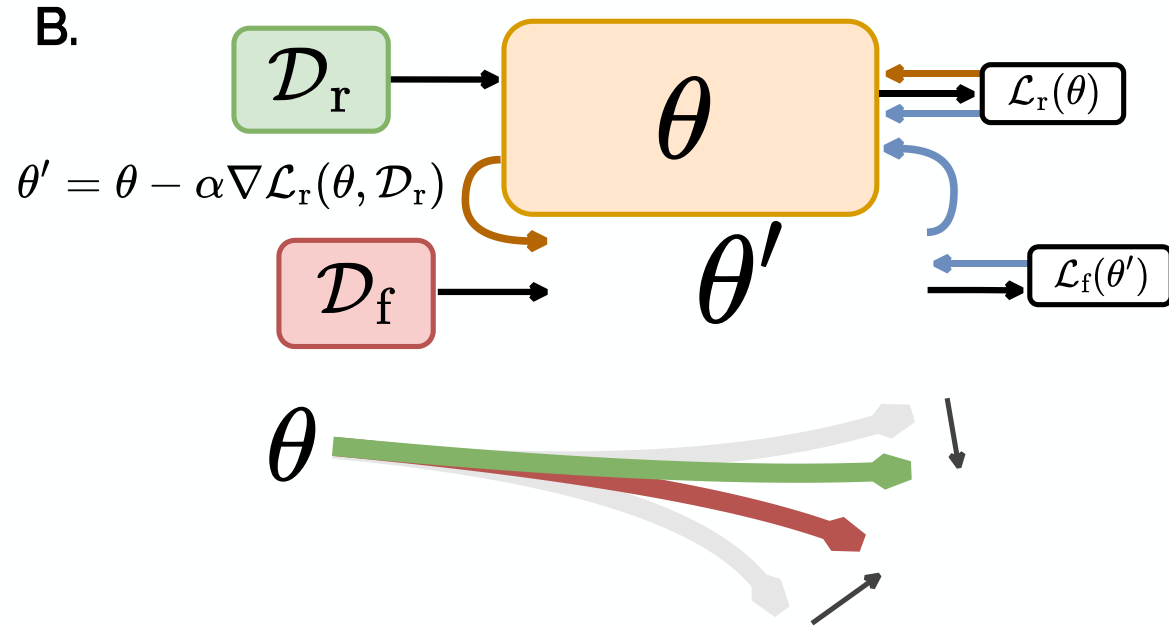
$$-\alpha \nabla (\nabla \mathcal{L}_f(\theta) \cdot \nabla \mathcal{L}_r(\theta))$$

- Encourages *maximizing* the inner product  $\langle \nabla \mathcal{L}_f(\theta) \cdot \nabla \mathcal{L}_r(\theta) \rangle$  reducing gradient conflicts.



# Our Proposition

## Learning to Unlearn while Retaining



# Quantitative Results

## Experiments: Classification & Generative Models

### Classification: Random data forgetting on CIFAR-10 and CIFAR-100

Table 1. Performance comparison of different MU methods for image classification under 10% (left) and 50% (right) *random data forgetting* scenarios on CIFAR-10 [31] (top) and CIFAR-100 [31] (bottom) using ResNet-18 [18]. Results are reported in the format  $a \pm b$ , where  $a$  denotes the mean and  $b$  represents the standard deviation over 10 independent trials. A smaller performance gap relative to Retrain indicates better MU method performance. The metric **Avg. Gap** quantifies this gap by computing the average absolute performance differences across the considered evaluation metrics (see Section 5). Best results highlighted in **Maroon** and second best in **Navy**.

Method	Random Data Forgetting (10%)					Random Data Forgetting (50%)				
	UA ( $\uparrow$ )	TA ( $\uparrow$ )	RA ( $\uparrow$ )	MIA ( $\uparrow$ )	Avg. Gap ( $\downarrow$ )	UA ( $\uparrow$ )	TA ( $\uparrow$ )	RA ( $\uparrow$ )	MIA ( $\uparrow$ )	Avg. Gap ( $\downarrow$ )
CIFAR 10										
Retrain	5.19 $\pm$ 0.53	94.26 $\pm$ 0.14	100.00 $\pm$ 0.00	13.05 $\pm$ 0.64	0	7.83 $\pm$ 0.26	91.71 $\pm$ 0.30	100.00 $\pm$ 0.00	19.13 $\pm$ 0.55	0
FT [63]	0.85 $\pm$ 0.46	93.83 $\pm$ 0.45	99.84 $\pm$ 0.11	3.01 $\pm$ 0.93	3.74	0.50 $\pm$ 0.33	94.32 $\pm$ 0.07	99.96 $\pm$ 0.03	2.31 $\pm$ 1.08	6.70
GA [59]	0.34 $\pm$ 0.23	94.57 $\pm$ 0.01	99.62 $\pm$ 0.25	0.91 $\pm$ 0.29	4.42	0.40 $\pm$ 0.27	94.55 $\pm$ 0.06	99.62 $\pm$ 0.26	0.96 $\pm$ 0.40	7.20
IU [30]	1.92 $\pm$ 2.10	91.91 $\pm$ 2.73	98.01 $\pm$ 2.26	4.01 $\pm$ 3.44	4.16	2.46 $\pm$ 1.99	91.10 $\pm$ 5.25	97.62 $\pm$ 1.98	5.25 $\pm$ 3.01	5.56
BE [5]	0.59 $\pm$ 0.38	93.79 $\pm$ 0.15	99.41 $\pm$ 0.38	16.16 $\pm$ 0.78	2.19	0.43 $\pm$ 0.28	94.28 $\pm$ 0.04	99.59 $\pm$ 0.28	10.82 $\pm$ 0.89	4.67
BS [5]	0.40 $\pm$ 0.25	94.24 $\pm$ 0.07	99.56 $\pm$ 0.54	4.46 $\pm$ 0.33	3.46	0.42 $\pm$ 0.28	94.44 $\pm$ 0.03	99.60 $\pm$ 0.27	1.99 $\pm$ 0.08	6.92
$\ell_1$ -sparse [39]	5.83 $\pm$ 0.49	90.64 $\pm$ 0.52	96.64 $\pm$ 0.54	11.87 $\pm$ 0.61	2.20	2.58 $\pm$ 0.60	92.10 $\pm$ 0.24	98.89 $\pm$ 0.15	6.59 $\pm$ 0.80	4.82
SalUn [10]	1.93 $\pm$ 0.42	93.92 $\pm$ 0.25	99.89 $\pm$ 0.07	17.93 $\pm$ 0.37	2.15	7.85 $\pm$ 1.18	88.15 $\pm$ 0.90	95.02 $\pm$ 0.98	19.30 $\pm$ 2.81	<b>2.18</b>
SHs [65]	4.60 $\pm$ 1.48	92.92 $\pm$ 0.48	98.93 $\pm$ 0.57	9.56 $\pm$ 2.13	<b>1.62</b>	7.98 $\pm$ 5.31	88.32 $\pm$ 4.24	94.00 $\pm$ 4.87	15.52 $\pm$ 6.43	3.29
LUR (Ours)	5.52 $\pm$ 2.16	92.95 $\pm$ 0.29	99.21 $\pm$ 0.27	11.93 $\pm$ 1.01	<b>0.89</b>	6.79 $\pm$ 0.81	90.23 $\pm$ 0.63	97.19 $\pm$ 0.72	13.98 $\pm$ 0.63	<b>2.62</b>
CIFAR 100										
Retrain	24.87 $\pm$ 0.85	74.69 $\pm$ 0.08	99.98 $\pm$ 0.01	50.22 $\pm$ 0.62	0	32.83 $\pm$ 0.14	67.27 $\pm$ 0.45	99.99 $\pm$ 0.01	60.76 $\pm$ 0.21	0
FT [63]	2.02 $\pm$ 1.36	75.28 $\pm$ 0.12	99.95 $\pm$ 0.02	9.64 $\pm$ 3.60	16.01	1.83 $\pm$ 1.20	75.36 $\pm$ 0.36	99.97 $\pm$ 0.01	9.26 $\pm$ 2.84	22.65
GA [59]	2.00 $\pm$ 1.34	75.59 $\pm$ 0.11	98.24 $\pm$ 1.16	5.00 $\pm$ 2.25	17.68	1.85 $\pm$ 1.23	75.50 $\pm$ 0.10	98.22 $\pm$ 1.17	4.94 $\pm$ 1.96	24.2
IU [30]	4.33 $\pm$ 4.82	72.13 $\pm$ 4.58	96.14 $\pm$ 4.51	9.43 $\pm$ 5.98	16.93	3.14 $\pm$ 2.19	72.08 $\pm$ 2.41	97.17 $\pm$ 2.00	8.20 $\pm$ 4.10	22.47
BE [5]	2.06 $\pm$ 1.38	74.16 $\pm$ 0.09	98.12 $\pm$ 1.24	7.60 $\pm$ 3.05	16.96	2.65 $\pm$ 1.60	67.84 $\pm$ 0.58	97.27 $\pm$ 1.62	8.62 $\pm$ 2.19	21.40
BS [5]	2.35 $\pm$ 1.48	73.20 $\pm$ 0.18	97.93 $\pm$ 1.30	8.24 $\pm$ 3.23	17.01	4.69 $\pm$ 1.47	68.12 $\pm$ 0.18	95.41 $\pm$ 1.46	10.07 $\pm$ 1.99	21.07
$\ell_1$ -sparse [39]	3.65 $\pm$ 0.67	70.06 $\pm$ 0.46	96.35 $\pm$ 0.67	21.33 $\pm$ 1.95	14.59	9.83 $\pm$ 2.43	69.73 $\pm$ 1.27	97.35 $\pm$ 0.89	21.72 $\pm$ 1.44	16.79
SalUn [10]	11.44 $\pm$ 1.18	71.34 $\pm$ 0.48	99.40 $\pm$ 0.35	74.66 $\pm$ 2.48	10.45	15.19 $\pm$ 0.91	64.94 $\pm$ 0.48	98.89 $\pm$ 0.48	73.86 $\pm$ 1.98	<b>8.54</b>
SHs [65]	31.24 $\pm$ 1.81	73.17 $\pm$ 0.24	99.24 $\pm$ 0.30	42.42 $\pm$ 2.06	<b>4.11</b>	20.27 $\pm$ 2.28	67.58 $\pm$ 1.76	84.64 $\pm$ 2.79	28.68 $\pm$ 2.53	15.08
LUR (Ours)	29.57 $\pm$ 0.26	73.02 $\pm$ 0.18	99.29 $\pm$ 0.06	41.44 $\pm$ 0.10	<b>3.96</b>	32.68 $\pm$ 1.75	63.02 $\pm$ 0.90	87.18 $\pm$ 0.74	45.69 $\pm$ 2.79	<b>8.07</b>



# Quantitative Results

## Experiments: Classification & Generative Models

- Classification Class-wise forgetting on CelebA-HQ-FIR (face recognition benchmark)

Table 2. Performance comparison of different MU methods for image classification under *class-wise data forgetting* on Celeb-HQ-FIR [33, 47] using ResNet-34 [18]. The content in follow the same format of Table 1. Best results highlighted in **Maroon** and second best in **Navy**.

Method	Random Class (Identity) Forgetting (10%)					Random Class (Identity) Forgetting (50%)				
	UA ( $\uparrow$ )	TA ( $\uparrow$ )	RA ( $\uparrow$ )	MIA ( $\uparrow$ )	Avg. Gap ( $\downarrow$ )	UA ( $\uparrow$ )	TA ( $\uparrow$ )	RA ( $\uparrow$ )	MIA ( $\uparrow$ )	Avg. Gap ( $\downarrow$ )
Retrain	100.00 $\pm$ 0.00	87.02 $\pm$ 0.80	99.96 $\pm$ 0.01	100.00 $\pm$ 0.00	0	100.00 $\pm$ 0.00	88.09 $\pm$ 1.37	99.98 $\pm$ 0.03	100.00 $\pm$ 0.00	0
FT [63]	0.06 $\pm$ 0.12	88.59 $\pm$ 0.59	99.97 $\pm$ 7.02	5.28 $\pm$ 2.03	49.06	0.02 $\pm$ 0.03	90.71 $\pm$ 1.27	99.98 $\pm$ 0.03	3.08 $\pm$ 0.24	49.46
GA [59]	12.4 $\pm$ 8.71	81.22 $\pm$ 2.11	99.74 $\pm$ 0.26	51.37 $\pm$ 5.96	35.56	0.04 $\pm$ 0.02	88.41 $\pm$ 0.40	99.98 $\pm$ 0.03	2.44 $\pm$ 0.43	49.46
IU [30]	11.08 $\pm$ 10.25	70.24 $\pm$ 11.77	95.27 $\pm$ 5.07	29.59 $\pm$ 18.59	45.20	9.63 $\pm$ 8.78	68.40 $\pm$ 7.91	94.80 $\pm$ 6.61	30.10 $\pm$ 9.65	46.29
BE [5]	30.93 $\pm$ 2.73	44.11 $\pm$ 2.08	95.58 $\pm$ 1.23	46.24 $\pm$ 5.90	42.53	0.06 $\pm$ 0.02	83.12 $\pm$ 1.68	99.97 $\pm$ 0.02	3.62 $\pm$ 0.52	50.33
BS [5]	1.82 $\pm$ 1.92	81.92 $\pm$ 0.27	99.86 $\pm$ 0.03	45.93 $\pm$ 5.11	39.36	0.02 $\pm$ 0.03	87.80 $\pm$ 0.95	99.98 $\pm$ 0.03	2.76 $\pm$ 0.35	49.38
$\ell_1$ -sparse [39]	1.19 $\pm$ 0.72	89.37 $\pm$ 0.70	99.97 $\pm$ 0.00	76.78 $\pm$ 5.66	31.10	23.86 $\pm$ 3.63	90.29 $\pm$ 1.05	99.92 $\pm$ 0.10	99.86 $\pm$ 0.19	19.64
SalUn [10]	100.00 $\pm$ 0.00	78.36 $\pm$ 1.34	96.90 $\pm$ 1.11	100.00 $\pm$ 0.00	2.93	45.10 $\pm$ 2.60	90.92 $\pm$ 1.66	99.98 $\pm$ 0.03	99.95 $\pm$ 0.00	14.45
SHs [65]	98.48 $\pm$ 2.73	80.18 $\pm$ 6.60	97.20 $\pm$ 3.81	99.83 $\pm$ 0.35	<b>2.82</b>	99.24 $\pm$ 0.52	81.64 $\pm$ 3.75	99.14 $\pm$ 0.95	100.00 $\pm$ 0.00	<b>2.01</b>
<b>LUR (Ours)</b>	100.00 $\pm$ 0.00	86.61 $\pm$ 1.01	99.97 $\pm$ 0.00	100.00 $\pm$ 0.00	<b>0.10</b>	99.75 $\pm$ 0.20	91.64 $\pm$ 0.74	99.97 $\pm$ 0.02	100.00 $\pm$ 0.00	<b>0.95</b>



# Quantitative Results

## Experiments: Classification & Generative Models

- Generative Models: Diffusion (Stable Diffusion)

Table 4. Class-wise forgetting performance on Imagenette [26] using SD [50]. Best results highlighted in **Maroon** and second best in **Navy**.

Forget Class	ESD [12]		FMN [68]		SalUn [10]		LUR (Ours)	
	UA (↑)	FID (↓)	UA (↑)	FID (↓)	UA (↑)	FID (↓)	UA (↑)	FID (↓)
Tench	99.40	1.22	42.40	1.63	100.00	2.53	100.00	0.74
English Springer	100.00	1.02	27.20	1.75	100.00	0.79	100.00	0.97
Cassette Player	100.00	1.84	93.80	0.80	99.80	0.91	99.80	0.99
Chain Saw	96.80	1.48	48.40	0.94	100.00	1.58	100.00	1.30
Church	98.60	1.91	23.80	1.32	99.60	0.90	100.00	1.04
French Horn	99.80	1.08	45.00	0.99	100.00	0.94	100.00	0.75
Garbage Truck	100.00	2.71	41.40	0.92	100.00	0.91	100.00	0.94
Gas Pump	100.00	1.99	53.60	1.30	100.00	1.05	100.00	0.88
Golf Ball	99.60	0.80	15.40	1.05	98.80	1.45	100.00	0.88
Parachute	99.80	0.91	34.40	2.33	100.00	1.16	99.80	1.29
Average	99.40	1.50	42.54	1.30	<b>99.82</b>	<b>1.22</b>	<b>99.96</b>	<b>0.98</b>



# Qualitative Results



Figure 4. *Quantitative and qualitative evaluation on I2P [53] benchmark.* **A.** Evaluation of the amount of nudity content detected by the NudeNet classifier [1] for each unlearning method. The bars represent the percentage decrease in the number of images from each nudity class compared to SD [50]. **B.** Generated images from SD with and without MU. Unlearning methods: SalUn [10], SHs [65], and **LUR** (Ours). Each column shows images from different MU methods using the same prompt ( $P_i$ ) and seed. Prompt details in Appendix Table 5 (Supplementary Material).



# Qualitative Results

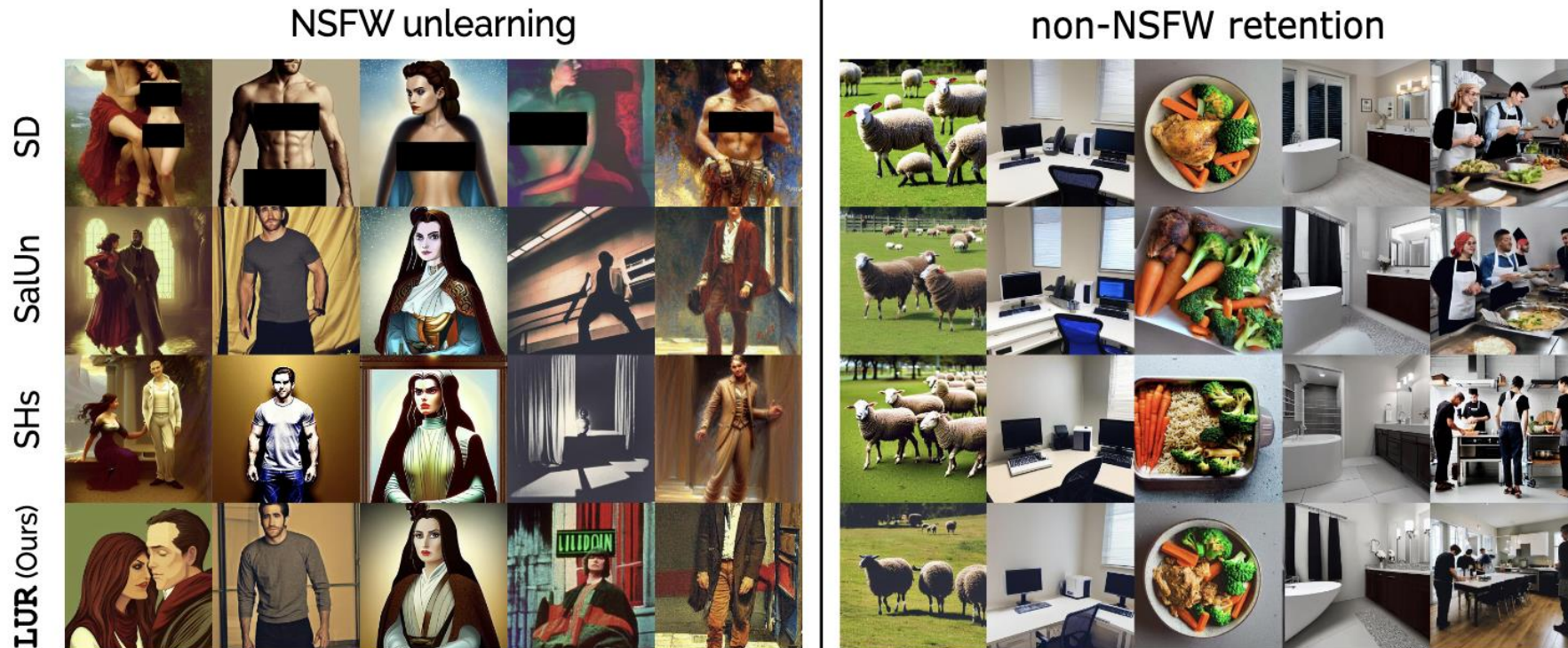


Figure 8. Example generations from prompts in I2P [53] (left) and COCO-10k [37, 69] (right) after unlearning the concept of *nudity*. Each column represents the generation from one prompts with a fixed seed. The prompts corresponding to the generated images are provided in Table 10.



# Conclusion & Future Directions

- **Main Contributions:**

- **LUR** framework aligns conflicting gradients for retain vs. forget → fosters a more robust approximate unlearning.
- **Implicit Gradient Product Regularization** elegantly emerges from the two-step update.
- Demonstrated **wide applicability** (discriminative and generative tasks), consistently improving over baseline and current MU methods.

- **Looking Ahead:**

- Potential for multi-concept unlearning, language model unlearning, and other multimodal settings (VLMs).



# ***Thank You***

