

# Sa1Un: Empowering Machine Unlearning via Gradient-based Weight Saliency in Both Image Classification and Generation

## Spotlight



Paper

Chongyu Fan<sup>1,\*</sup>, Jiancheng Liu<sup>1,\*</sup>, Yihua Zhang<sup>1</sup>,  
Eric Wong<sup>2</sup>, Dennis Wei<sup>3</sup>, Sijia Liu<sup>1,3</sup>

<sup>1</sup>Michigan State University, <sup>2</sup>University of Pennsylvania, <sup>3</sup>IBM Research

\*Equal contribution



Code

*{fanchong2, liujia45}@msu.edu*

# What is Machine Unlearning (MU)?

- Eliminate undesirable data influence (e.g., sensitive or illegal information) and associated model capabilities, while maintaining utility.
- Applications: Removing sensitive data information, copyright protection, harmful content degeneration, etc.



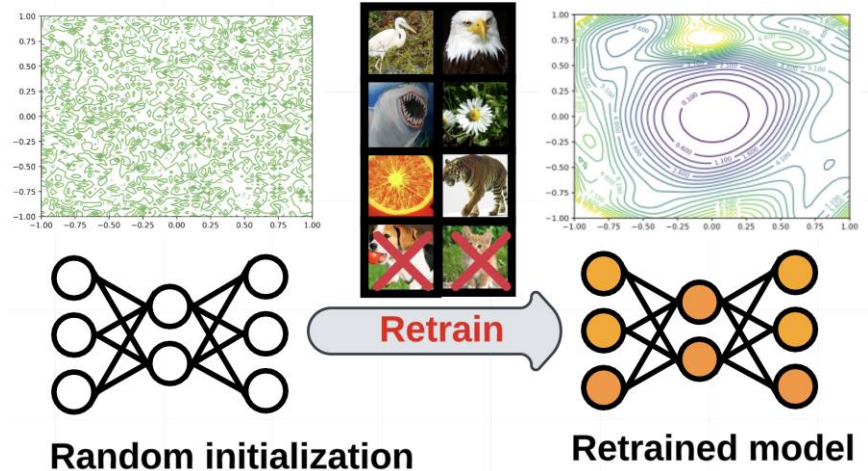
Generation examples by Stable Diffusion pre/after applying MU (Sa1Un).

**Forgetting Objective:** (Left) Concept "*Nudity*"; (Mid) Object "*Dog*"; (Right) Style "*Sketch*".



# Why and Why Not Retrain?

- Retrain model from scratch over retaining dataset (after removing data to be unlearned) is considered as an **optimal** MU method.
- Limitation: Lacks training efficiency, particularly for large-scale deep models

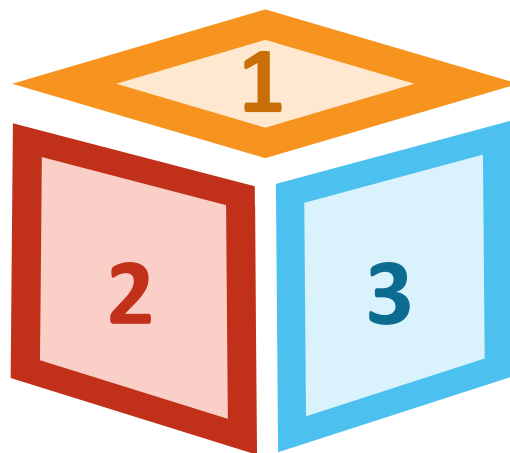


# How to Define the “GOOD” in MU?

**Computation  
Efficiency**

**Generalization  
Fidelity**

Can unlearned models  
still generalize?



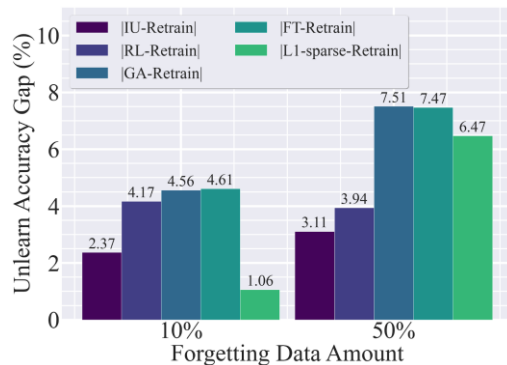
**Unlearning Efficacy**

Is the impact of forgetting  
data points truly removed?



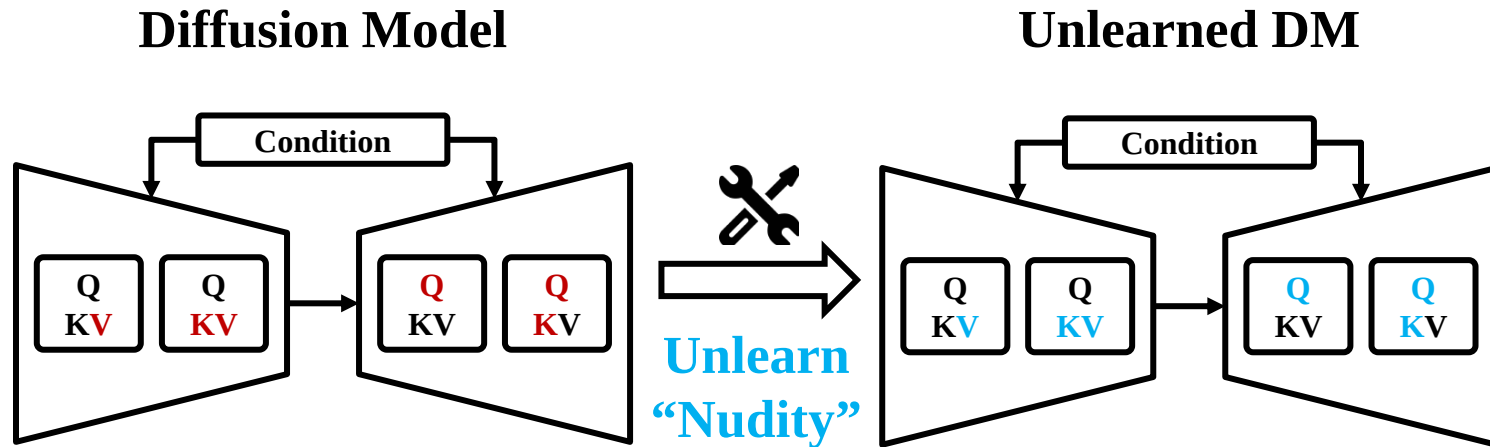
# Limitations of Current MU Methods

- Existing MU methods lack generality in harder tasks.
  - Large forgetting ratio
  - Image generation tasks
- Hard to strike the balancing point between unlearning and generalization.
  - Tend to either **over-forget** (e.g., GA , RL)
  - Or **under-forget** (e.g., FT , l1-sparse)



	Original	Retrain	GA	RL	FT
Forgetting class: "airplane"					

# Sa1Un: Weight Saliency Is the Key to MU



# Sa1Un: Weight Saliency Is the Key to MU

- Use the gradient of the forgetting loss with respect to the model weights  $\theta$  under the forgetting dataset.
- Apply a hard thresholding to obtain the weight saliency mask.
- By fixing the low-saliency parameters, achieving an accurate unlearning.

$$\mathbf{m}_S = \mathbb{1} \left( \left| \nabla_{\theta} l_f(\theta; \mathcal{D}_f) \Big|_{\theta=\theta_o} \right| \geq \gamma \right)$$
$$\theta_u = \underbrace{\mathbf{m}_S \odot (\Delta\theta + \theta_o)}_{\text{salient weights}} + \underbrace{(\mathbf{1} - \mathbf{m}_S) \odot \theta_o}_{\text{original weights}}$$



# Sa1Un: Weight Saliency Is the Key to MU

- Integrate weight saliency with random labeling (RL) provides a promising MU solution both in image classification and image generation.
- Classification:
  - Sa1Un assigns a random label to each forgetting data point and then fine-tunes the salient weights on the randomly relabeled dataset.

$$\underset{\Delta\theta}{\text{minimize}} L_{\text{Sa1Un}}^{(1)}(\theta_u) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_f, y' \neq y} [\ell_{\text{CE}}(\theta_u; \mathbf{x}, y')] + \alpha \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_r} [\ell_{\text{CE}}(\theta_u; \mathbf{x}, y)]$$





# Sa1Un: Weight Saliency Is the Key to MU

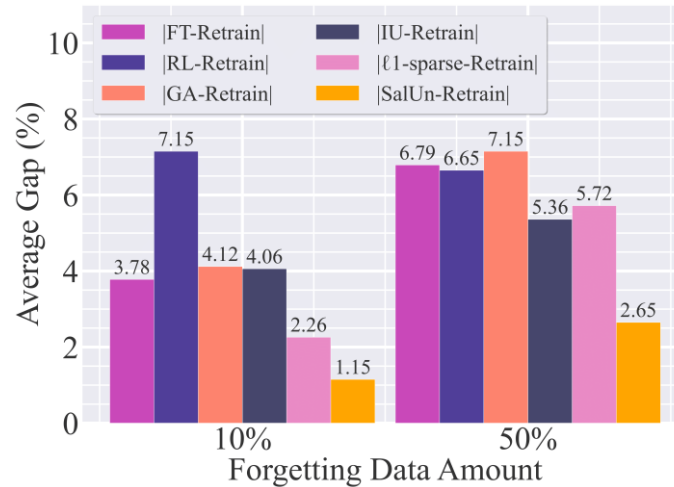
- Integrate weight saliency with random labeling (RL) provides a promising MU solution both in image classification and image generation.
- Generation:
  - Sa1Un associates each image  $x$  in forgetting concept  $c$  with a misaligned concept  $c'$ .

$$\underset{\Delta\theta}{\text{minimize}} L_{\text{Sa1Un}}^{(2)}(\theta_u) := \mathbb{E}_{(\mathbf{x}, c) \sim \mathcal{D}_f, t, \epsilon \sim \mathcal{N}(0,1), c' \neq c} [\|\epsilon_{\theta_u}(\mathbf{x}_t | c') - \epsilon_{\theta_u}(\mathbf{x}_t | c)\|_2^2] + \beta \ell_{\text{MSE}}(\theta_u; \mathcal{D}_r)$$



# Sa1Un in Image Classification

- Sa1Un demonstrates consistent forgetting performance across different forgetting data amounts.



# Sa1Un in Image Generation

- Class-wise forgetting performance in image generation



Example of forgetting "Church" class.

Forget. Class	Sa1Un		ESD		FMN	
	UA (↑)	FID (↓)	UA (↑)	FID (↓)	UA (↑)	FID (↓)
Tench	<b>100.00</b>	2.53	99.40	<b>1.22</b>	42.40	1.63
English Springer	<b>100.00</b>	<b>0.79</b>	100.00	1.02	27.20	1.75
Cassette Player	99.80	0.91	<b>100.00</b>	1.84	93.80	<b>0.80</b>
Chain Saw	<b>100.00</b>	1.58	96.80	1.48	48.40	<b>0.94</b>
Church	<b>99.60</b>	<b>0.90</b>	98.60	1.91	23.80	1.32
French Horn	<b>100.00</b>	<b>0.94</b>	99.80	1.08	45.00	0.99
Garbage Truck	<b>100.00</b>	<b>0.91</b>	100.00	2.71	41.40	0.92
Gas Pump	<b>100.00</b>	<b>1.05</b>	100.00	1.99	53.60	1.30
Golf Ball	98.80	1.45	<b>99.60</b>	<b>0.80</b>	15.40	1.05
Parachute	<b>100.00</b>	1.16	99.80	<b>0.91</b>	34.40	2.33
Average	<b>99.82</b>	<b>1.22</b>	99.40	1.49	42.54	1.30

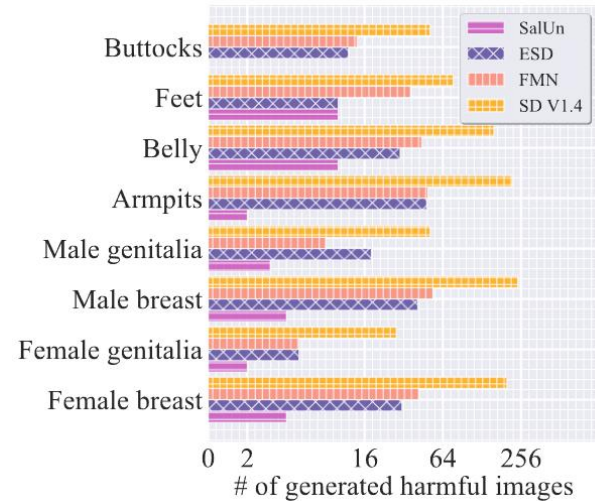


# SaUn in Image Generation

- Concept forgetting performance in image generation



Example of forgetting "Nudity" concept.



# Summary

- Weight saliency helps Machine Unlearning
- Limitations
  - Limited scope
  - Fine-grained masking methods



Met dank  
obrigada

terima kasih

multumesc

ありがとう

谢谢

ngiyabonga suksema

baie dankie

molte grazie

Thank

merci

감사합니다

obrigado

Danke schön!

謝謝

You

Благодарность

شكرًا

gracias

Спасиби

Dziękuję

dank u

mahalo

tusind tak