



# Targeted Unlearning with Single Layer Unlearning Gradient

Zikui Cai<sup>1,2 \*</sup>, Yaoteng Tan<sup>1 \*</sup>, M. Salman Asif<sup>1</sup>

<sup>1</sup> University of California, Riverside

<sup>2</sup> University of Maryland

*\* Equal Contribution*



# Challenges of Machine Unlearning

- Removing the influence of a specific subset of training data (forget set), while retaining overall model utility on the retain set.

$$\min_{\theta} \underbrace{\frac{1}{N_r} \sum_{(x_r, y_r) \in D_r} \ell(F_{\theta}(x_r), y_r)}_{\mathcal{L}_{\text{retain}}} - \underbrace{\frac{\alpha}{N_f} \sum_{(x_f, y_f) \in D_f} \ell(F_{\theta}(x_f), y_f)}_{\mathcal{L}_{\text{forget}}}$$

Fine-tuning (FT)      Gradient ascent (GA)

Two stage: Gradient ascent + Fine-tuning (GAFT)

- Challenges:
  - **High Computational Cost:** involves iterative gradient calculation and model update over the whole model
  - **Side Effects:** Removing one concept can lead to degraded performance on unrelated tasks due to interconnected representations
  - **Robustness Issues:** Unlearned concepts can be recovered via careful probing or attacks

# Our Three Main Objectives

## Computational Efficiency

Minimize computational overhead compared to full retraining:

- Achieve unlearning with minimal parameter updates
- Enable practical unlearning for large-scale models

## Effective Unlearning

Ensure complete removal of targeted concepts:

- Achieve zero forget accuracy on target concepts
- Prevent information leakage through related concepts
- Resist recovery through adversarial techniques

## Targeted Removal with Minimal Side Effects

Maintain model utility while removing specific content:

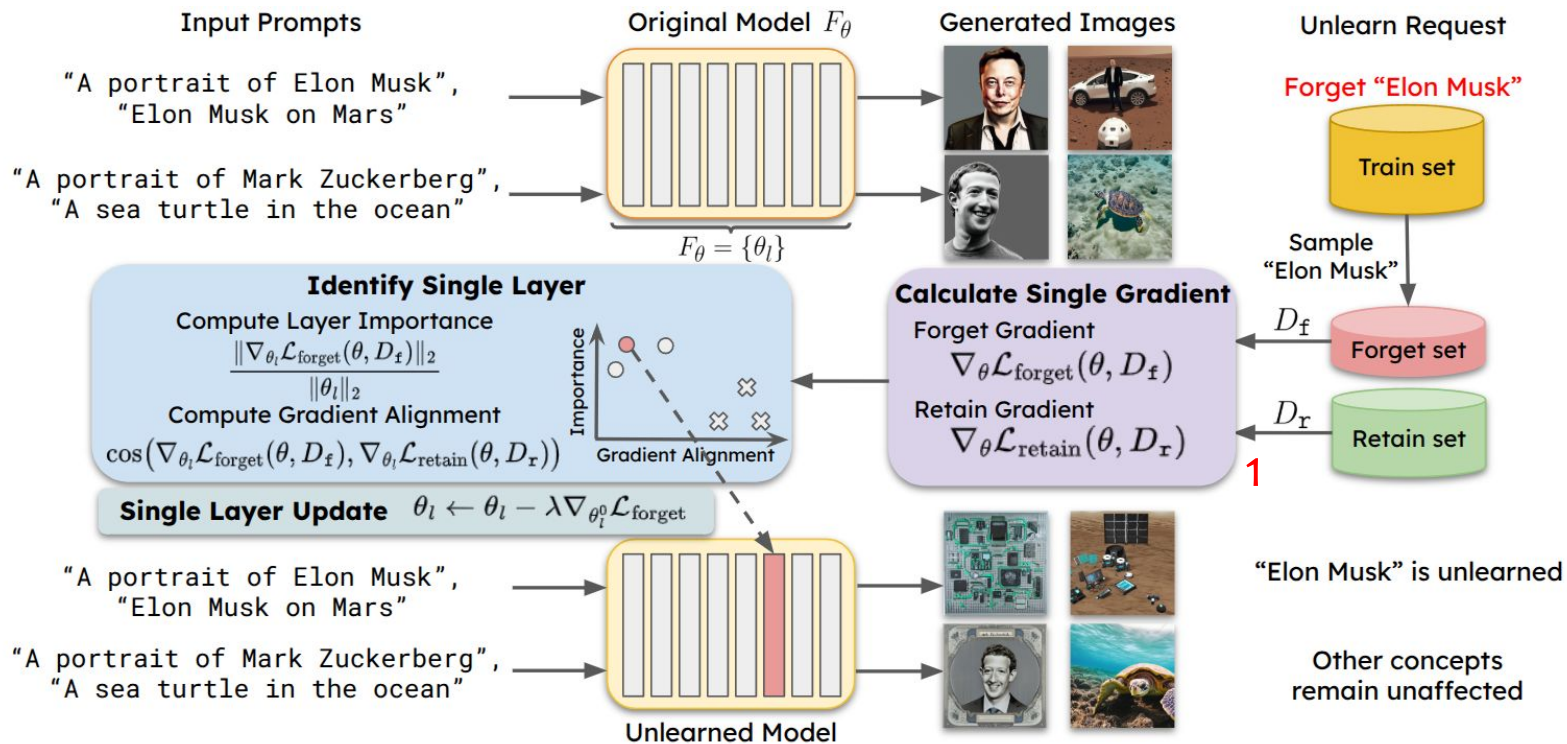
- Preserve performance on unrelated tasks
- Maintain accuracy on retain set
- Ensure precision in targeting only unwanted concepts

# Introduction

- Goal:
  - Balance unlearning effectiveness & retaining model utility
  - Increase efficiency (computational time + resource)
- Questions
  - Can we select the most critical model part(s) for unlearning?
  - Can we do update that part efficiently?
- Proposed method:
  - Achieves unlearning with reduced model parameter manipulation
  - Requires one-time gradient calculation and a single-layer update.

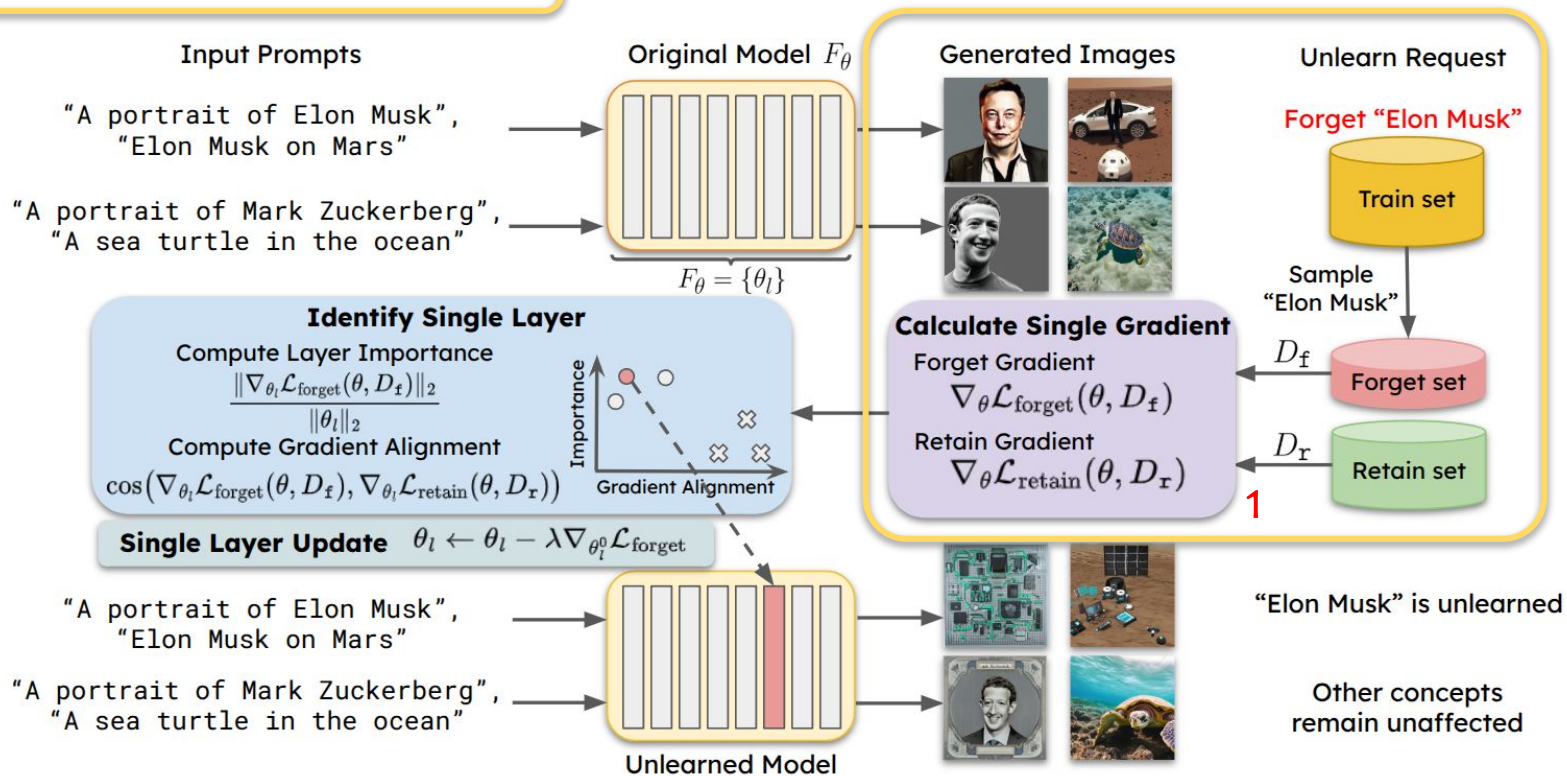
# Framework: Single Layer Unlearning Gradient (SLUG 🐌)

- 1 Compute gradients →



# Framework: Single Layer Unlearning Gradient (SLUG 🐌)

- 1 Compute gradients →



# Compute Gradient

- Sample the forget and retain sets
  - Forget set:  
Images related to unlearning requests (e.g., “Elon Musk”)
  - Retain set:  
Subset of training set excluding the forget data

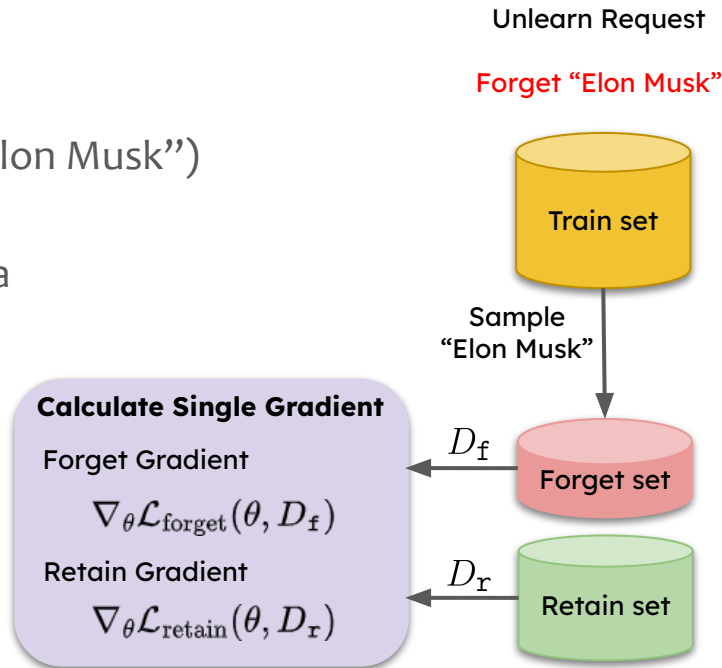
- Compute forget and retain gradients

- Forget loss:  
Alignments of text and image embedding

$$\mathcal{L}_{\text{forget}}(\mathbf{v}_i, \mathbf{t}_j) = 1 - \cos(\mathbf{v}_i, \mathbf{t}_j)$$

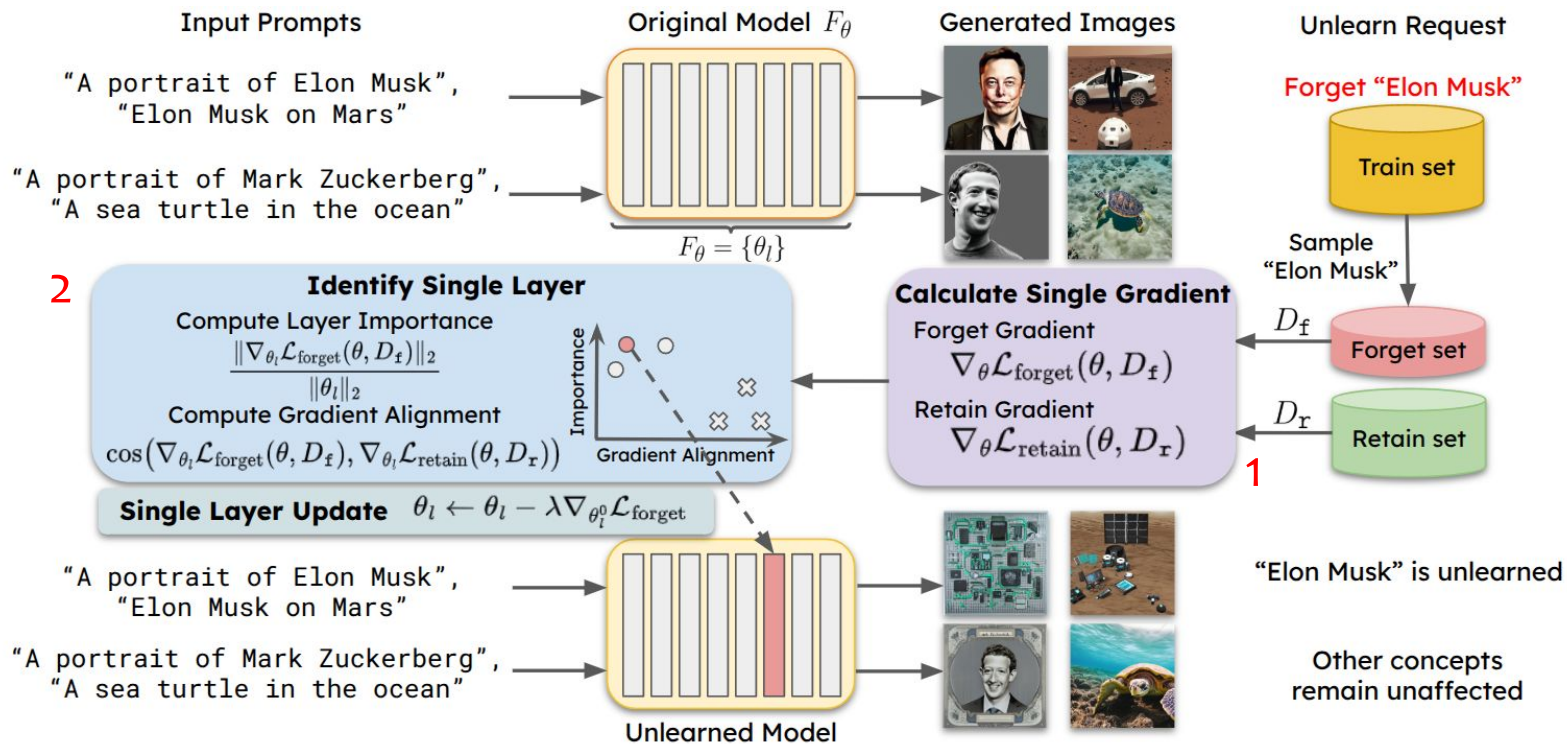
- Retain loss:  
Image and Text contrastive loss (CLIP training)

$$\mathcal{L}_{\text{retain}} = \frac{1}{2N} \sum_{i=1}^N (\ell_{i2t}(i) + \ell_{t2i}(i))$$



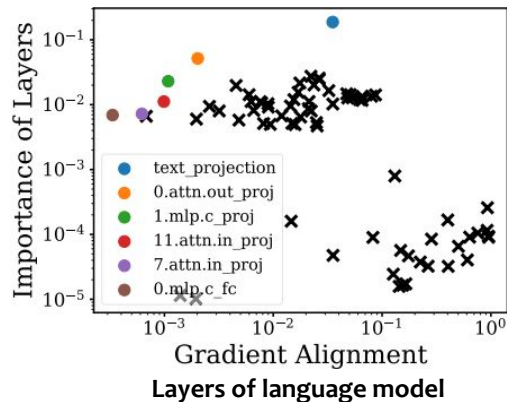
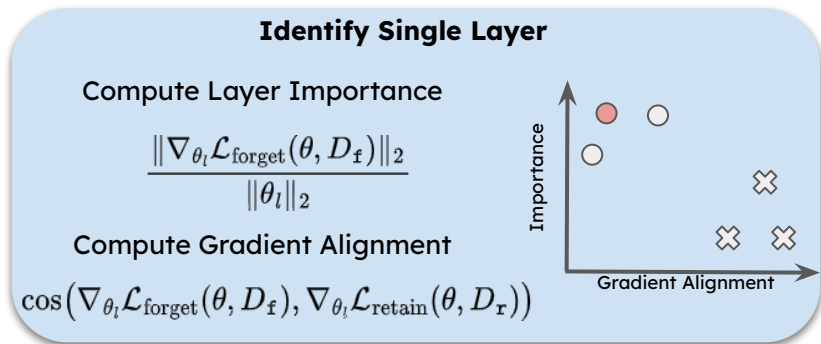
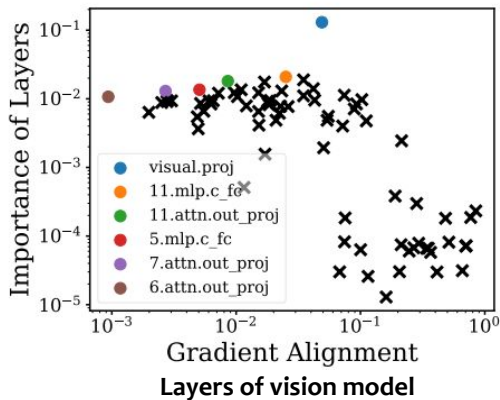
# Framework: Single Layer Unlearning Gradient (SLUG )

- 1 Compute gradients  $\rightarrow$  2 Identify layer to update



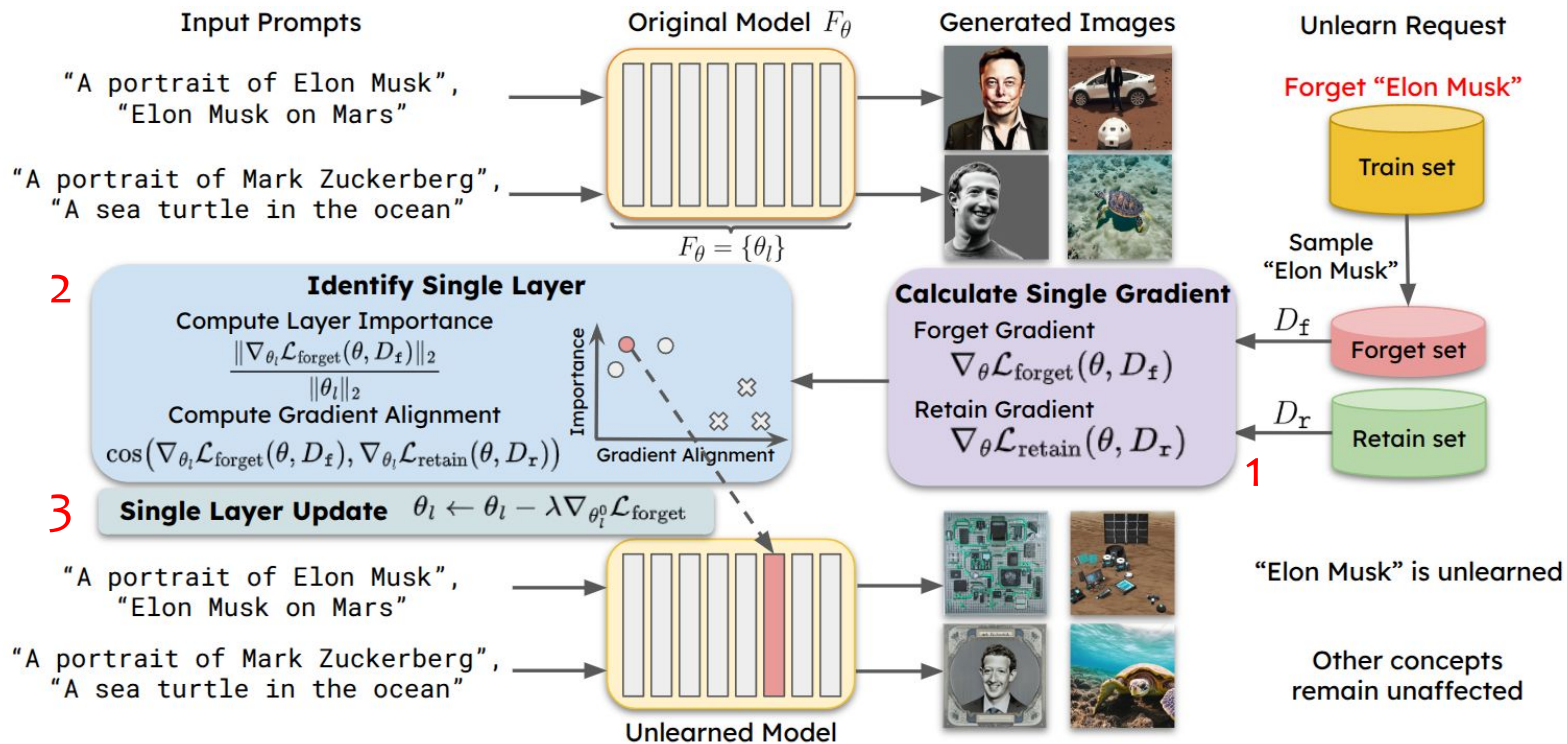
# Layer Identification

- Identify the most critical part to update
- Utilize Forget and Retain gradients
  - **Layer importance:**  
Higher the more impactful on unlearning
  - **Gradient alignment:**  
Lower the least impactful on retaining
  - **Pareto-front:** Layers that are well-balanced for unlearning and retaining



# Framework: Single Layer Unlearning Gradient (SLUG )

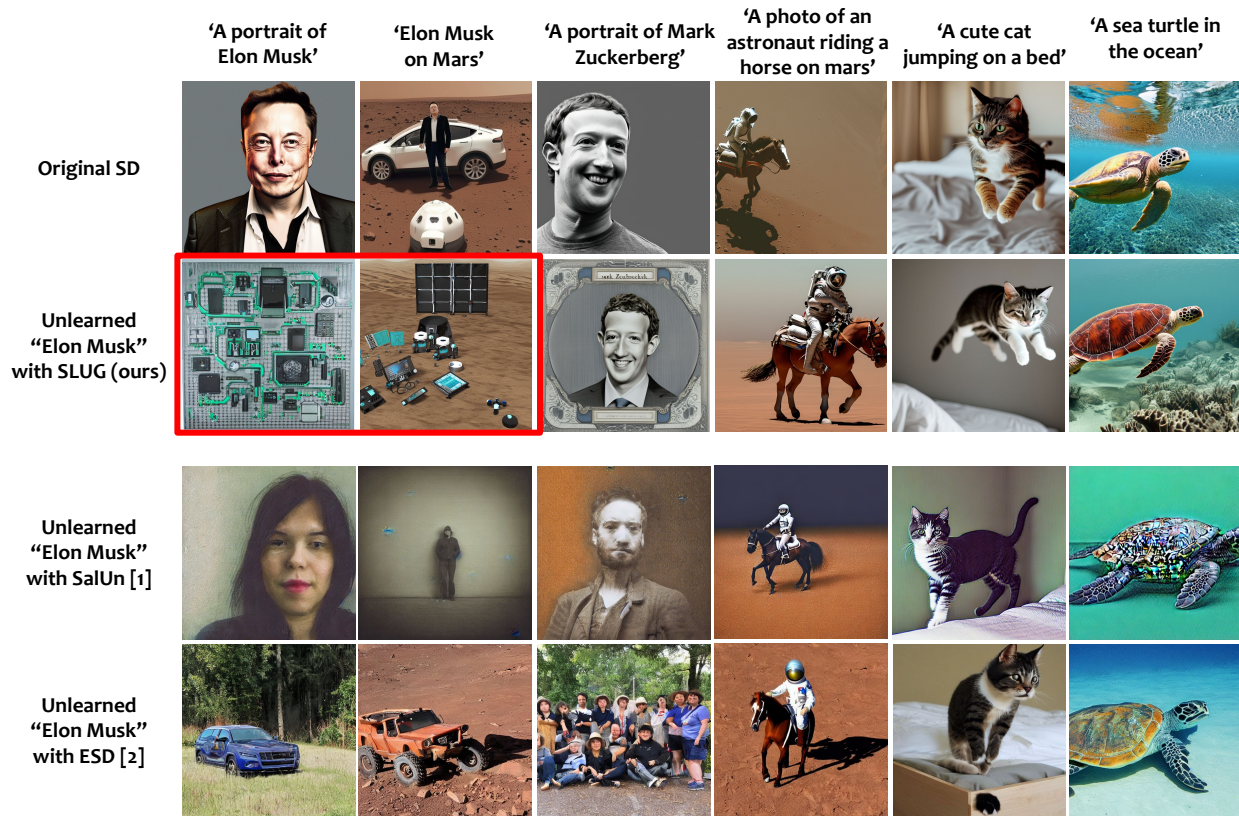
- 1 Compute gradients  $\rightarrow$  2 Identify layer to update  $\rightarrow$  3 Search proper step-size



# Unlearning Text-to-Image Models

## ● Setup

- Unlearning identity “Elon Musk” on Stable Diffusion
- Comparing methods
  - SalUn [1]
  - ESD [2]



Examples of Stable Diffusion unlearned “Elon Musk”

[1] SalUn (Fan et. al, ICLR 2024)

[2] ESD (Gandikota et. al, ICCV 2023)

# More examples

- Setup

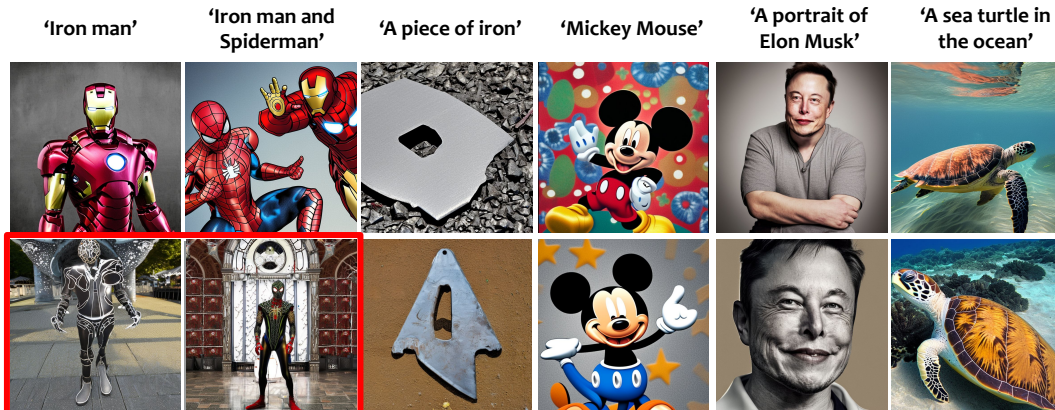
- Unlearning intellectual properties that have copyright on Stable Diffusion

- “Iron man”\*

- “Mickey Mouse”\*

Original SD

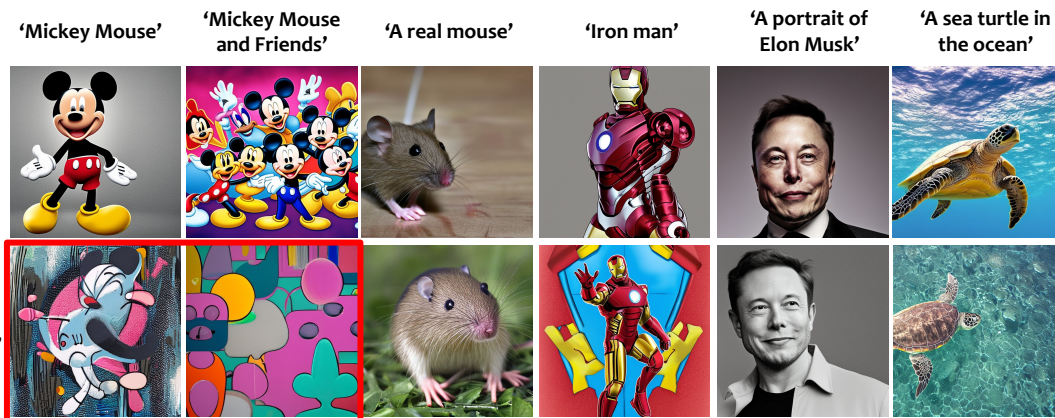
Unlearned  
“Iron man”



Examples of Stable Diffusion unlearned “Iron Man”\*

Original SD

Unlearned  
“Mickey Mouse”



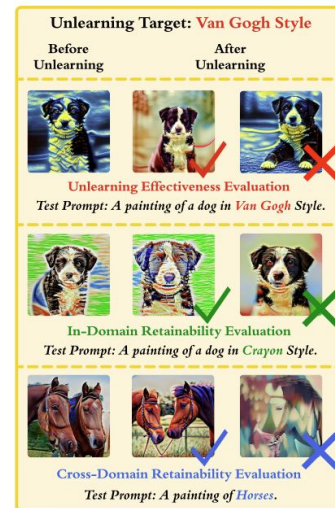
Examples of Stable Diffusion unlearned “Mickey Mouse”\*

\*Iron Man: Marvel Comics character

\*Mickey Mouse: Walt Disney character

# Results: UnlearnCanvas benchmark

- Setup
  - UA (unlearning acc)
  - IRA (In domain retain acc)
  - CRA (cross domain retain acc)
- Main takeaway
  - SLUG achieve the best trade-off between efficiency and effectiveness




Quantitative evaluation of SLUG unlearning on UnlearnCanvas benchmark

| Method      | Effectiveness    |         |        |                   |         |        |         | Efficiency   |                 |                  |
|-------------|------------------|---------|--------|-------------------|---------|--------|---------|--------------|-----------------|------------------|
|             | Style Unlearning |         |        | Object Unlearning |         |        | FID (↓) | Time (s) (↓) | Memory (GB) (↓) | Storage (GB) (↓) |
| UA (↑)      | IRA (↑)          | CRA (↑) | UA (↑) | IRA (↑)           | CRA (↑) |        |         |              |                 |                  |
| ESD [9]     | 98.58%           | 80.97%  | 93.96% | 92.15%            | 55.78%  | 44.23% | 65.55   | 6163         | 17.8            | 4.3              |
| FMN [40]    | 88.48%           | 56.77%  | 46.60% | 45.64%            | 90.63%  | 73.46% | 131.37  | 350          | 17.9            | 4.2              |
| UCE [10]    | 98.40%           | 60.22%  | 47.71% | 94.31%            | 39.35%  | 34.67% | 182.01  | 434          | 5.1             | 1.7              |
| CA [18]     | 60.82%           | 96.01%  | 92.70% | 46.67%            | 90.11%  | 81.97% | 54.21   | 734          | 10.1            | 4.2              |
| SalUn [7]   | 86.26%           | 90.39%  | 95.08% | 86.91%            | 96.35%  | 99.59% | 61.05   | 667          | 30.8            | 4.0              |
| SEOT [21]   | 56.90%           | 94.68%  | 84.31% | 23.25%            | 95.57%  | 82.71% | 62.38   | 95           | 7.34            | 0.0              |
| SPM [25]    | 60.94%           | 92.39%  | 84.33% | 71.25%            | 90.79%  | 81.65% | 59.79   | 29700        | 6.9             | 0.0              |
| EDiff [36]  | 92.42%           | 73.91%  | 98.93% | 86.67%            | 94.03%  | 48.48% | 81.42   | 1567         | 27.8            | 4.0              |
| SHS [35]    | 95.84%           | 80.42%  | 43.27% | 80.73%            | 81.15%  | 67.99% | 119.34  | 1223         | 31.2            | 4.0              |
| SLUG (Ours) | 86.29%           | 84.59%  | 88.43% | 75.43%            | 77.50%  | 81.18% | 75.97   | 39           | 3.61            | 0.04             |

# Unlearning Image-to-Text Models

## ● Setup

- Unlearning identities on VLMs (LLaVA 1.5-7B)
  - “Elon Musk”
  - “Taylor Swift”




Input Image

User Prompt: “What’s the name of the person in this image?”

Answer of LLaVA Pretrained: “The person in this image is Elon Musk.”


Answer of LLaVA Unlearned: “The person in this image is Michael Jackson.”



User Prompt: “What’s the name of the person in this image?”

Answer of LLaVA Pretrained: “The person in this image is Mark Zuckerberg.”

Answer of LLaVA Unlearned: “The person in this image is Mark Zuckerberg.”



User Prompt: “What’s the content of the image?”


Answer of LLaVA Pretrained: “The image features a stop sign on a street corner, with a car driving by in the background. There are also people walking on the sidewalk, and a lion statue nearby. The scene takes place in an Asian city, with a red and”

Answer of LLaVA Unlearned: “The image features a stop sign on a street corner, with a red building in the background. There is also a statue of a lion nearby.”

Examples of LLaVA unlearned “Elon Musk”

Quantitative evaluation of unlearned LLaVA 1.5 models

| Model                    | Forget Accuracy (%) (↓) | VLM Benchmark Score (↑) |                |       |              |
|--------------------------|-------------------------|-------------------------|----------------|-------|--------------|
|                          |                         | MME Cognition           | MME Perception | GQA   | MMBench (en) |
| Original LLaVA-1.5       | 99.50                   | 323.57                  | 1481.21        | 61.28 | 62.97        |
| Unlearned “Elon Musk”    | <b>3.0</b>              | 298.57                  | 1354.61        | 60.70 | 61.34        |
| Unlearned “Taylor Swift” | <b>2.0</b>              | 334.64                  | 1336.09        | 60.72 | 60.14        |
| Average                  | 2.5                     | 316.61                  | 1345.35        | 60.71 | 60.74        |




Input Image

User Prompt: “What’s the name of the person in this image?”

Answer of LLaVA Pretrained: “The person in this image is Taylor Swift.”


Answer of LLaVA Unlearned: “The person in this image is a woman.”



User Prompt: “What’s the name of the person in this image?”

Answer of LLaVA Pretrained: “The person in this image is Queen Elizabeth II.”


Answer of LLaVA Unlearned: “The person in this image is Queen Elizabeth II.”



User Prompt: “What’s the name of the person in this image?”

Answer of LLaVA Pretrained: “The person in this image is Marilyn Monroe.”


Answer of LLaVA Unlearned: “The person in this image is Marilyn Monroe.”



User Prompt: “What’s the content of the image?”

Answer of LLaVA Pretrained: “The image is a cartoon depiction of a baseball game, featuring various characters and elements from the sport.”

Answer of LLaVA Unlearned: “The image is a cartoon or comic strip featuring a baseball game, with various characters and animals participating in the game.”



User Prompt: “What is unusual about this image?”

Answer of LLaVA Pretrained: “The unusual aspect of this image is that the man is reading a book while sitting in a swimming pool. It is not common to see someone engaging in such an activity, as people typically read books in more comfortable and dry environments”

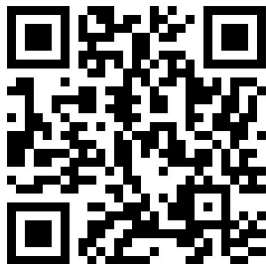
Answer of LLaVA Unlearned: “The unusual aspect of this image is that a man is reading a book while sitting underwater. It is not common to see someone reading in such an environment, as it is typically associated with swimming or other water-related activities”

Examples of LLaVA unlearned “Taylor Swift”

# Conclusion

- Summary
  - We propose a single layer, one-step unlearning for vision-language foundation models that significantly improves the unlearning efficiency
  - Our framework is scalable to different vision-language tasks
- More details

Paper



Code



Project page

