# Model Sparsity Can Simplify Machine Unlearning

Jinghan Jia, Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, Pranay Sharma, Sijia Liu

37th Conference on Neural Information Processing Systems (2023)

# TABLE OF CONTENTS

**01**

Introduction
- Overview
- Problem

**02**

Related Work
- Unlearning
- Pruning

**03**

Proposal
- Challenges
- Solutions

**04**

Method
- Theory
- Paradigms

**05**

Experiments
- Setup
- Results

**06**

Conclusion
- Limitations
- Future Work

# 01

## Introduction

# 1.1 OVERVIEW

### What Is Machine Unlearning (MU)?

- Reverse learning to remove influence of specific examples from an already trained model

### Why Is MU Needed?

- Recent data regulation requirements e.g., privacy of data A.K.A. "the right to be forgotten", corrupted data, etc.

### Why Is Retraining Not Enough?

- Direct and optimal but unreasonable computational costs
- Approximate and fast but effective methods required

# 1.2 PROBLEM

*Driving Question (Q)*

**(Q)** Is there a theoretically-grounded and broadly-applicable method to improve approximate unlearning across different unlearning criteria?

# 1.2 PROBLEM

## *Problem Setup*

- Training dataset of $N$ points (with labels):  $\mathcal{D} = \{x_i\}_{i=1}^{N}$

- **F**orgetting dataset (to be scrubbed):  $\mathcal{D}_f \subseteq \mathcal{D}$

  ↳  $\mathcal{D}_f$ can be *class-wise* or *random data* forgetting

- **R**emaining dataset (for new model):  $\mathcal{D}_r = \mathcal{D} \setminus \mathcal{D}_f$

- **O**riginal model parameters (on $\mathcal{D}$):  $\theta_o$

- **U**nlearned model parameters (on $\mathcal{D}_r$):  $\theta_u$

## *Problem Definition*

- Generate $\theta_u$ from $\theta_o$ accurately and efficiently

# 02

# Related Work

## 2.1 Unlearning

- Approximate MU
- Evaluation of MU
- Other Paradigms

## 2.2 Pruning

- Performance Impact
- Pruning in MU

# 2.1 UNLEARNING

***Approximate MU Methods***

- **Fine-tuning (FT):** fine-tunes $\theta_o$ on $\mathcal{D}_r$ for a few epochs to obtain $\theta_u$ and causes "catastrophic forgetting" as in continual learning

- **Gradient ascent (GA):** reverses training on $\mathcal{D}_f$ by adding the gradient back to $\theta_o$ to increase loss on $\mathcal{D}_f$ and obtain $\theta_u$

- **Fisher forgetting (FF):** adds Gaussian noise to $\theta_o$ to perturb dependency on $\mathcal{D}_f$ based on Fisher Information Matrix to obtain $\theta_u$

- **Influence unlearning (IU):** leverages influence function to find change in $\theta_o$ if $\mathcal{D}_f$ is removed and subtracts $\mathcal{D}_f$ impact to obtain $\theta_u$

# 2.1 UNLEARNING

- **More on influence unlearning (IU):** relates to an important line of research in MU ($\varepsilon - \delta$ forgetting); defined as:

$$\Delta(\mathbf{w}) \coloneqq \theta(\mathbf{w}) - \theta_o \approx \mathrm{H}^{-1}\nabla_\theta L(1/N - \mathbf{w}, \theta_o) \qquad \textit{update of } \theta_o \textit{ to } \theta(\mathbf{w})$$

where: $\quad \theta(\mathbf{w}) = \mathrm{argmin}_\theta L(\mathbf{w}, \theta) \qquad\qquad \textit{weighted ERM training}$

$$L(\mathbf{w}, \theta) = \sum_{i=1}^{N}[w_i \ell_i(\theta, z_i)]$$

$$w_i \in [0, 1], \quad \mathbf{1}^T\mathbf{w} = 1 \qquad\qquad \textit{influence of } z_i; \textit{ normal}$$

$$\mathrm{H}^{-1} = \left(\nabla_{\theta,\theta}^2 L(1/N, \theta_o)\right)^{-1} \qquad\qquad \textit{inv-Hessian; expensive} \longrightarrow$$

$$\xrightarrow{\text{scrub } \mathcal{D}_f} \theta_u = \theta_o + \Delta(\mathbf{w}_{MU}), \qquad \mathbf{w}_{MU} \in [0, 1]^N, \qquad w_{MU,i} = \mathbb{I}_{D_r}(i)/|D_r|$$

*current authors use WoodFisher approx. for implementation*

# 2.1 UNLEARNING

**Full-Stack MU Evaluation**

*Evaluated compared to Retrain gold-standard metrics*

- **Unlearning accuracy (UA):** defined as $\mathrm{UA}(\theta_u) = 1 - \mathrm{Acc}_{\mathrm{D_f}}(\theta_u)$

- **Membership inference attack (MIA-Efficacy):** confidence-based MIA predictor against $\theta_u$ on $\mathcal{D}_f$; defined as $\mathrm{T}N/|D_f|$

- **Remaining accuracy (RA):** accuracy of $\theta_u$ on $\mathcal{D}_r$; fidelity of MU

- **Testing accuracy (TA):** generalization of $\theta_u$ on dataset outside of $\mathcal{D}$; tested on all test data except in class-wise forgetting

- **Run-time efficiency (RTE):** computation efficacy of MU

# 2.1 UNLEARNING

***Other Paradigms***

- **Differential privacy:** protecting individuals in dataset; probabilistic

- **Federated learning:** training with decentralized edge devices

- **Graph neural networks:** processing data represented by graphs

- **Adversarial ML:** attacking ML algorithms for info or manipulation

- **Conditional generative models:** generating concepts to image

- **Understanding data influence:** influence function approach, defense against data poisoning, fair learning, transfer learning

# 2.2 PRUNING

- Necessary due to constraints on computation, memory, etc.

- Equates to sparsification A.K.A. weight sparsity/pruning

### *Performance Impact*

- Lottery ticket hypothesis (LTH) demonstrated the feasibility of co-improving test accuracy and efficiency (sparsity) of model

- Impact of pruning has been investigated in improving:
    *generalization, robustness, fairness, interpretability, model explanation, privacy, loss landscape*

- Privacy gains imply data influence connected to sparsification

*winning ticket is a sparse subnetwork with equal or better accuracy*

# 2.2 PRUNING

***Pruning in MU***

- A search for insights from pruning for unlearning

- **Wang et al. 2022:** removing channels of a DNN showed an unlearning benefit in federated learning.

- **Ye et al. 2022:** filter pruning was introduced in lifelong learning to detect "pruning identified exemplars" that are easy to forget

# 03

# Proposal

## 3.1  Challenges

- Works' Limitations
- General Challenges

## 3.2  Solutions

- Schematic Overview
- Summary of Work

# 3.1 CHALLENGES

### *Limitations of Related Work*

- **Exact unlearning (Retrain):** large computational overhead
- **Approximate unlearning:** speed and ease at the cost of efficacy

*FT and DP impractical against attacks, GA's effectiveness of unlearning can be improved, FF low parallel efficiency and dependent on parameters, IU requires model and training assumptions*

| Unlearning Methods | UA | MIA-Efficacy | RA | TA | RTE |
|---|---|---|---|---|---|
| FT | ✓ | | ✓ | ✓ | 0.06× |
| GA | ✓ | ✓ | ✓ | ✓ | 0.02× |
| FF | ✓ | | ✓ | ✓ | 0.9 × |
| IU | ✓ | | | ✓ | 0.08× |
| Ours | ✓ | ✓ | ✓ | ✓ | 0.07× |

*Table 1*

# 3.1 CHALLENGES

***General Challenges Categorized***

1. Performance of approximate unlearning can heavily rely on the configuration of algorithmic parameters.

    *e.g. FF regularization parameter for each data-model setup*

2. Effectiveness of scheme can vary significantly across different unlearning evaluation criteria, and their trade-offs are not well understood.

    *e.g. high efficacy neither implies nor precludes high fidelity*

# 3.2 SOLUTIONS

## *Schematic Overview of Proposal on Model-Sparsity Driven MU*

- Going beyond data-centric
- Use of model sparsity in MU to incorporate its impact
- Use of full stack evaluation metrics for a multi-faceted analysis of models
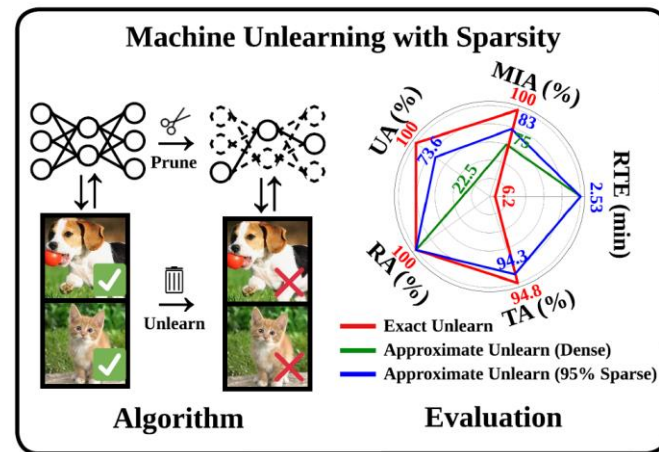- Comparison of various MU approaches to sparsity



Figure 1

# 3.2 SOLUTIONS

***Summary of Work***

1.  Systematically deriving theoretical connections of unlearning and pruning (rather than focusing on a specific application).

2.  Practically demonstrating the effects of the theory in closing the gap between approximate unlearning and exact unlearning.

3.  Developing a new paradigm termed "prune first, then unlearn" and a novel "sparsity-aware unlearning" framework to explore different methods of employing sparsification.

4.  Performing extensive experiments across diverse datasets, models, and unlearning scenarios.

# 04

## Method

### 4.1 Theory

- Fundamental Idea
- Error Analysis
- Expansion

### 4.2 Paradigms

- Prune First, then MU
- Sparsity-Aware MU
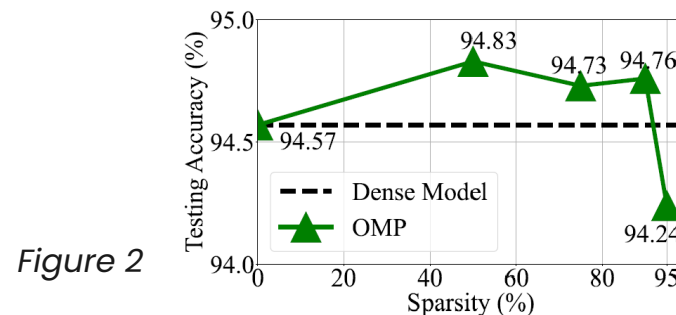
# 4.1 UNLEARNING

## *Fundamental Idea*

- Based on related work, we theorize sparsity boosts multi-criteria unlearning and closes approximation gap while being efficient
- First, we prove this is the case both theoretically and practically
- Then, we develop paradigms to employ this finding

**But how do we prove this?**

*Figure 2*

# 4.1 UNLEARNING

***Error Analysis to Prove Application of Sparsity***

- Let us use:
  - *Unrolling stochastic gradient descent to derive unlearning error given by weight difference in scrubbing a single data point*
  - *One-shot magnitude pruning to infuse model sparsity to SGD*
- Let us assume:
  - *Sparse pattern from OMP as binary mask :*      $\mathbf{m}$, $m_i \in [0,1]$
  - *Model parameters for every $m_i$:*      $\boldsymbol{\theta}$, $\theta_i$
  - *Sparse model with zeroed weights:*      $\mathbf{m} \odot \boldsymbol{\theta}$

# 4.1 UNLEARNING

***Error Analysis to Prove Application of Sparsity (cont.)***

- Using SGD, error between GA-unlearned and Retrain standard:

$$e(\mathbf{m}) = O\big(\eta^2 t \parallel \mathbf{m} \odot (\boldsymbol{\theta}_t - \boldsymbol{\theta}_0) \parallel_2 \sigma(\mathbf{m})\big) \qquad \textit{unlearning error}$$

$$\text{where:} \qquad \sigma(\mathbf{m}) \coloneqq \max_j\{\sigma_j\big(\nabla^2_{\theta,\theta}\ell\big), if\ m_j \neq 0\} \qquad \textit{largest singular value}$$

- Clearly, unlearning error decreases as sparsity increases, unlike previously being proportional to model distance
- But, number of active singular values decreases as sparsity increases, causing possible generalization decrease (TA can tell)

# 4.1 UNLEARNING

**Error Analysis to Prove Application of Sparsity (*proof overview*)**

$$\theta'_t \approx \theta'_0 - \eta \mathbf{m} \odot \sum_{i=1}^{t-1} \nabla_{\boldsymbol{\theta}} \ell(\theta'_0, \hat{\mathbf{z}}_i) + \mathbf{m} \odot \left( \sum_{i=1}^{t-1} f(i) \right),$$

$$f(i) = -\eta \nabla^2_{\boldsymbol{\theta},\boldsymbol{\theta}} \ell(\theta'_0, \hat{\mathbf{z}}_i) \left( -\eta \sum_{j=0}^{i-1} \mathbf{m} \odot \nabla_{\boldsymbol{\theta}} \ell(\theta'_0, \hat{\mathbf{z}}_j) + \sum_{j=0}^{i-1} (\mathbf{m} \odot f(j)) \right),$$

$$e(\mathbf{m}) = \|\mathbf{e_m}(\theta_0, \{\hat{\mathbf{z}}_i\}, t, \eta)\|_2 = \left\| \mathbf{m} \odot \left( \sum_{i=1}^{t-1} f(i) \right) \right\|_2$$

$$\approx \eta^2 \left\| \mathrm{diag}(\mathbf{m}) \sum_{i=1}^{t-1} \nabla^2_{\boldsymbol{\theta},\boldsymbol{\theta}} \ell(\theta'_0, \hat{\mathbf{z}}_i) \sum_{j=0}^{i-1} \mathbf{m} \odot \nabla_{\boldsymbol{\theta}} \ell(\theta'_0, \hat{\mathbf{z}}_j) \right\|_2$$

*through triangle inequality…*

$$\leq \eta^2 \sigma(\mathbf{m}) \|\mathbf{m} \odot (\theta_t - \theta_0)\|_2 \frac{1}{t} \frac{t-1}{2} t = \frac{\eta^2}{2} (t-1) \|\mathbf{m} \odot (\theta_t - \theta_0)\|_2 \sigma(\mathbf{m}),$$

# 4.1 UNLEARNING

***Expansion of Theory to Practice***

- Does the above benefit of model sparsification in MU apply to other approximate unlearning methods besides GA?

- Let us check the common metrics i.e. unlearning efficacy (UA and MIA-Efficacy), fidelity (RA), and generalization (TA):
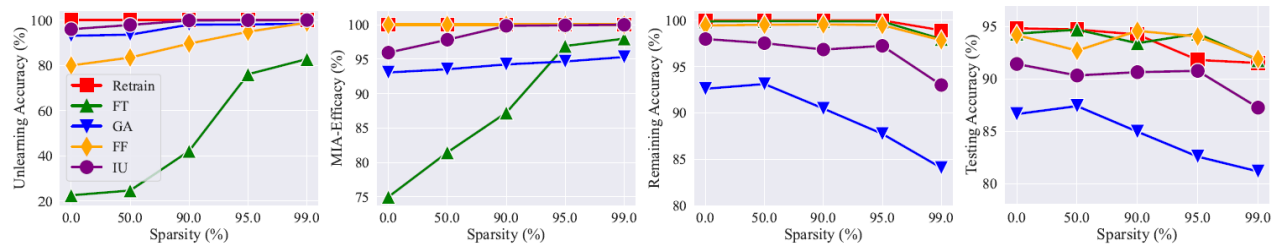


*Figure 3*

# 4.2 PARADIGMS

- Theoretically and with limited methods, sparsity worked

- But:

  ▪ **How does the choice of weight-pruning method impact the unlearning performance?**

  ▪ **Can sparsity-aware MU methods that directly scrub data influence from a dense model be developed?**

# 4.2 PARADIGMS

***Prune First, then Unlearn***                    (as we saw before)

- **What pruning to choose?**
  - Random initialization pruning before training?
  - Simultaneous pruning-training iterative magnitude pruning?

- **What matters?**
  1. Least dependence on $\mathcal{D}_f$
  2. Lossless generalization when pruning
  3. Pruning efficiency

# 4.2 PARADIGMS

***Prune First, then Unlearn (cont.)***

Choice of pruning methods based on criteria

- **SynFlow (synaptic flow pruning):** training-free pruning method at initialization, even without accessing the dataset (for ❶)
- **OMP (one-shot magnitude pruning):** performed over $\theta_o$ and depends on $\mathcal{D}_f$ but computationally light (for ❸) and has better generalization (for ❷)
- **IMP (iterative magnitude pruning):** not suitable for MU despite accuracy due to computation overhead and dependence on $\mathcal{D}_f$

# 4.2 PARADIGMS

## *Prune First, then Unlearn (cont.)*

- We compare the efficacy of FT-based MU on sparse models generated using different pruning methods (SynFlow, OMP, and IMP)
- We see:
  - IMP UA decrease due to reliance
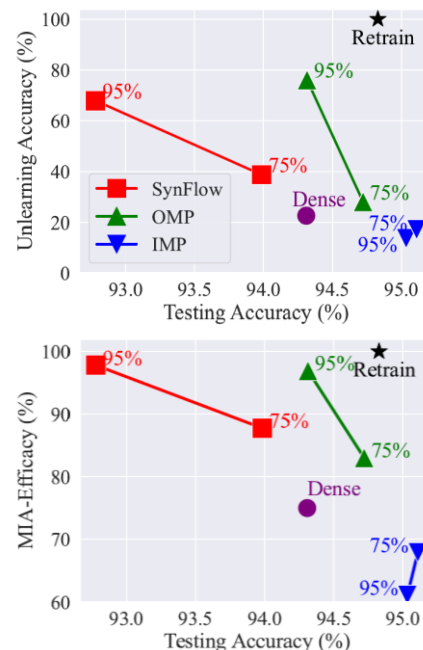  - OMP closer to Retrain
- OMP will be our choice due to balance



*Figure 4*

# 4.2 PARADIGMS

**Sparsity-Aware Unlearning**

- How about a method for simultaneous pruning and unlearning?
- Inspired by sparsity-inducing optimization, we integrate sparse penalty ($\ell_1$ norm-based) into unlearning objective function
- Therefore $\ell_1$**-sparse MU**:

$$\theta_u = \mathrm{argmin}_\theta L_u(\theta; \theta, D_r) + \gamma \parallel \theta \parallel_1$$

where:  $L_u(\theta; \theta, D_r)$  *FT objective function*

$\gamma > 0$  *regularization*

# 4.2 PARADIGMS

### *Sparsity-Aware Unlearning (cont.)*

- Choice of $\gamma$ matters and is a limitation; a spare-regularization scheduler can mitigate this issue with three schemes:

  1. Constant $\gamma$

  2. Linearly growing $\gamma$

  3. Linearly decaying $\gamma$                    *outperforms others!*

*Table 2*

| MU | UA | MIA-Efficacy | RA | TA | RTE (min) |
|---|---|---|---|---|---|
| Retrain | 5.41 | 13.12 | 100.00 | 94.42 | 42.15 |
| $\ell_1$-sparse MU + constant $\gamma$ | 6.60 (1.19) | 14.64 (1.52) | 96.51 (3.49) | 87.30 (7.12) | 2.53 |
| $\ell_1$-sparse MU + linear growing $\gamma$ | 3.80 (1.61) | 8.75 (4.37) | 97.13 (2.87) | 90.63 (3.79) | 2.53 |
| $\ell_1$-sparse MU + linear decaying $\gamma$ | **5.35 (0.06)** | **12.71 (0.41)** | **97.39 (2.61)** | **91.26 (3.16)** | 2.53 |

# 05

# Experiments

## 5.1 Setup

- Datasets & Models
- MU & Pruning
- Evaluation

## 5.2 Results

- Method Tests
- Applications
- Other Results

# 5.1 SETUP

### Datasets and Models

- Mainly image classification under CIFAR-10 using ResNet-18

- Other datasets and an alternate architecture also explored

### Unlearning and Pruning

- Class-wise and random data forgetting for FT, GA, FF, and IU

- Both paradigms "prune-first" and "sparse-aware" with OMP

### Evaluation

- UA, MIA-Efficacy, RA, TA, and RTE

- Whole proximity to Retrain gauged using Disparity Average

# 5.2 RESULTS

**Experiments on Proposed Methods**

- **Prune first, then unlearn:**

  *Performance gap reduces with sparsity but Retrain has 3% TA drop; FT & UI preserve TA with tradeoff; FF loss in random-data forgetting*

| MU | UA | | MIA-Efficacy | | RA | | TA | | Disparity Ave. ↓ | | RTE |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | DENSE | 95% Sparsity | DENSE | 95% Sparsity | DENSE | 95% Sparsity | DENSE | 95% Sparsity | DENSE | 95% Sparsity | (min) |
| Class-wise forgetting | | | | | | | | | | | |
| Retrain | $100.00_{\pm0.00}$ | $100.00_{\pm0.00}$ | $100.00_{\pm0.00}$ | $100.00_{\pm0.00}$ | $100.00_{\pm0.00}$ | $99.99_{\pm0.01}$ | $94.83_{\pm0.11}$ | $91.80_{\pm0.89}$ | 0.00 | 0.00 | 43.23 |
| FT | $22.53_{\pm8.16}$ (77.47) | $73.64_{\pm9.46}$ (26.36) | $75.00_{\pm14.68}$ (25.00) | $83.02_{\pm16.33}$ (16.98) | $99.87_{\pm0.04}$ (0.13) | $99.87_{\pm0.05}$ (0.12) | $94.31_{\pm0.19}$ (0.52) | $94.32_{\pm0.12}$ (2.52) | 25.78 | 11.50 | 2.52 |
| GA | $93.08_{\pm2.29}$ (6.92) | $98.09_{\pm1.11}$ (1.91) | $94.03_{\pm3.27}$ (5.97) | $97.74_{\pm2.24}$ (2.26) | $92.60_{\pm0.25}$ (7.40) | $87.74_{\pm0.27}$ (12.25) | $86.64_{\pm0.28}$ (8.19) | $82.58_{\pm0.27}$ (9.22) | 7.12 | 6.41 | 0.33 |
| FF | $79.93_{\pm8.92}$ (20.07) | $94.83_{\pm4.29}$ (5.17) | $100.00_{\pm0.00}$ (0.00) | $100.00_{\pm0.00}$ (0.00) | $99.45_{\pm0.24}$ (0.55) | $99.48_{\pm0.33}$ (0.51) | $94.18_{\pm0.08}$ (0.65) | $94.04_{\pm0.10}$ (2.24) | 5.32 | 1.98 | 38.91 |
| IU | $87.82_{\pm2.15}$ (12.18) | $99.47_{\pm0.15}$ (0.53) | $95.96_{\pm0.21}$ (4.04) | $99.93_{\pm0.04}$ (0.07) | $97.98_{\pm0.21}$ (2.02) | $97.24_{\pm0.13}$ (2.75) | $91.42_{\pm0.21}$ (3.41) | $90.76_{\pm0.18}$ (1.04) | 5.41 | 1.10 | 3.25 |
| Random data forgetting | | | | | | | | | | | |
| Retrain | $5.41_{\pm0.11}$ | $6.77_{\pm0.23}$ | $13.12_{\pm0.14}$ | $14.17_{\pm0.18}$ | $100.00_{\pm0.00}$ | $100.00_{\pm0.00}$ | $94.42_{\pm0.09}$ | $93.33_{\pm0.12}$ | 0.00 | 0.00 | 42.15 |
| FT | $6.83_{\pm0.51}$ (1.42) | $5.97_{\pm0.57}$ (0.80) | $14.97_{\pm0.62}$ (1.85) | $13.36_{\pm0.59}$ (0.81) | $96.61_{\pm0.25}$ (3.39) | $96.99_{\pm0.31}$ (3.01) | $90.13_{\pm0.26}$ (4.29) | $90.29_{\pm0.31}$ (3.04) | 2.74 | 1.92 | 2.33 |
| GA | $7.54_{\pm0.29}$ (2.13) | $5.62_{\pm0.46}$ (1.15) | $10.04_{\pm0.31}$ (3.08) | $11.76_{\pm0.52}$ (2.41) | $93.31_{\pm0.04}$ (6.69) | $95.44_{\pm0.11}$ (4.56) | $89.28_{\pm0.07}$ (5.14) | $89.26_{\pm0.15}$ (4.07) | 4.26 | 3.05 | 0.31 |
| FF | $7.84_{\pm0.71}$ (2.43) | $8.16_{\pm0.67}$ (1.39) | $9.52_{\pm0.43}$ (3.60) | $10.80_{\pm0.37}$ (3.37) | $92.05_{\pm0.16}$ (7.95) | $92.29_{\pm0.24}$ (7.71) | $88.10_{\pm0.19}$ (6.32) | $87.79_{\pm0.23}$ (5.54) | 5.08 | 4.50 | 38.24 |
| IU | $2.03_{\pm0.43}$ (3.38) | $6.51_{\pm0.52}$ (0.26) | $5.07_{\pm0.74}$ (8.05) | $11.93_{\pm0.68}$ (2.24) | $98.26_{\pm0.29}$ (1.74) | $94.94_{\pm0.31}$ (5.06) | $91.33_{\pm0.22}$ (3.09) | $88.74_{\pm0.42}$ (4.59) | 4.07 | 3.08 | 3.22 |

*Table 3*

# 5.2 RESULTS

**Experiments on Proposed Methods (cont.)**

- **Sparsity-aware unlearning:**

  Only comparing to FT due to previous results

*It outperforms FT in efficacy*

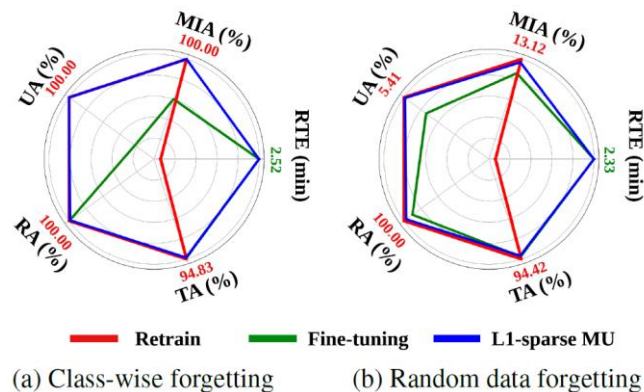*Closes gap with Retrain*

*Faces no computation loss*



*Figure 5*

# 5.2 RESULTS

***Uses In Different Applications***

- **MU for Trojan model cleanse:**

    Removing influence of poisoned backdoor data, manipulated by injecting trigger; can cause incorrect prediction with trigger

*FT decreases ASR*
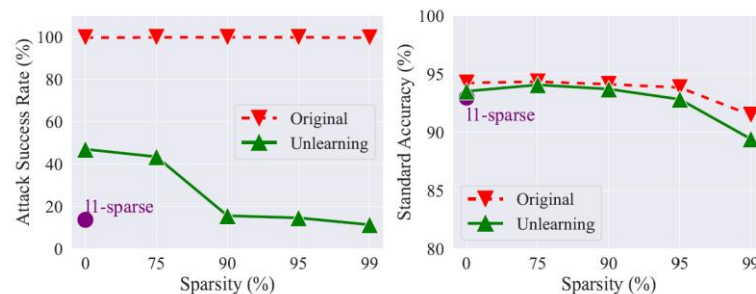
*FT has little SA loss*

$\ell_1$-*sparse performs similarly*

*Figure 6*

# 5.2 RESULTS

### Uses In Different Applications (cont.)

- **MU for transfer learning:**

  Mitigating impact of harmful data classes to enhance model's accuracy on other datasets after finetuning; we use $\ell_1$-*sparse*

  *Keeps accuracy of Retrain*

  *Has 2× speed up*

  *Suitable for large-scale*

| Forgetting class # | 0 Acc | 100 Acc | 100 Time | 200 Acc | 200 Time | 300 Acc | 300 Time |
|---|---|---|---|---|---|---|---|
| **OxfordPets** | | | | | | | |
| Method [51] | 85.70 | 85.79 | 71.84 | 86.10 | 61.53 | 86.32 | 54.53 |
| $\ell_1$-sparse MU | | 85.83 | 35.47 | 86.12 | 30.19 | 86.26 | 26.49 |
| **SUN397** | | | | | | | |
| Method [51] | 46.55 | 46.97 | 73.26 | 47.14 | 61.43 | 47.31 | 55.24 |
| $\ell_1$-sparse MU | | 47.20 | 36.69 | 47.25 | 30.96 | 47.37 | 27.12 |

*Table 4*

# 5.2 RESULTS

***Additional Results***

- **Model sparsity for data privacy:**

    Assessing MIA-Privacy to check how much gets leaked in MIA about $\mathcal{D}_r$; lower is better

    *Sparsity increase causes MIA-P decrease*

    *Approx. bests Retrain ($\mathcal{D}_r$ dependent)*
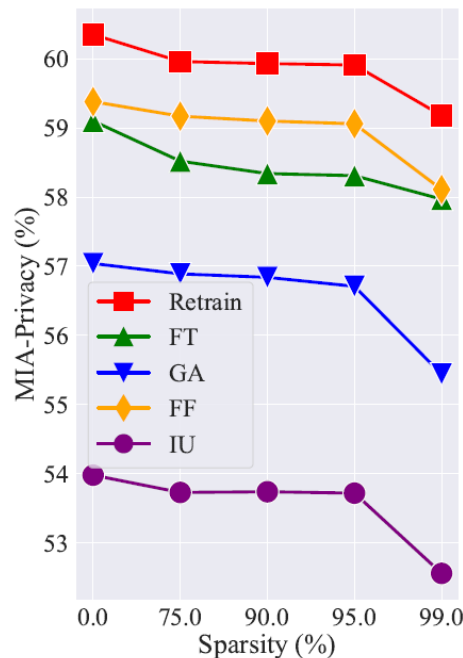
    *IU and GA are more private*



*Figure A1*

# 06

**Conclusion**

6.1  Limitations

6.2  Future Work

# 6.1 LIMITATIONS

- In the theoretical analysis of sparsity, it is unclear if it is better than other methods that improve performance, and if yes, why.

- $\varepsilon - \delta$ forgetting is briefly addressed but despite relevance, is not explored and its contribution is unevaluated.

- Current simulations are all on computer vision tasks which limits the application domain (though architectures are varied).

- Theoretical analysis is not performed on $\ell_1$-*sparse* update rule and might be difficult (similar to influence function methods relying on strongly convex loss).

# 6.2 FUTURE WORK

- The study indicates model modularity traits such as weight sparsity that could amplify MU and should be investigated.
- The application of this approach in other types of datasets and model architectures (e.g. language models) might face new challenges and reveal interesting results.

# THANKS!

**Do you have any questions?**

**Presentation by Maryam Rezaee**

Machine Learning Seminar | *Fall 1403*

Dr. Fatemeh SeyyedSalehi

*Sharif University of Technology*