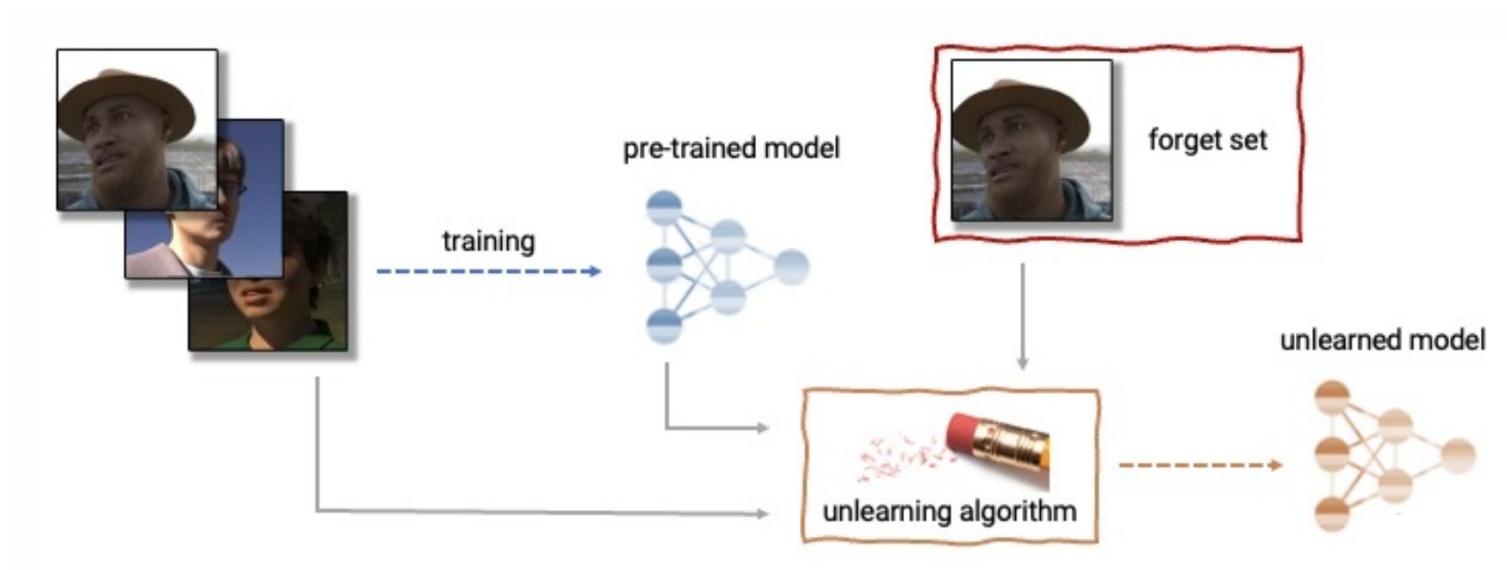


One-Shot Machine Unlearning with Mnemonic Code

ACML 2024

Tomoya Yamashita, Masanori Yamada, Takashi Shibata
NTT Laboratories

What is Machine Unlearning?



Challenges in Machine Unlearning

- We usually do not have access to the data.
- + we should not have access to the data.
- Computational costs

All above are addressed in this paper.

Fisher Information Matrix / Hessian Matrix

$$I(\theta) = \mathbb{E} \left[\left(\frac{\partial \log L(\theta)}{\partial \theta} \right) \left(\frac{\partial \log L(\theta)}{\partial \theta} \right)^T \right]$$

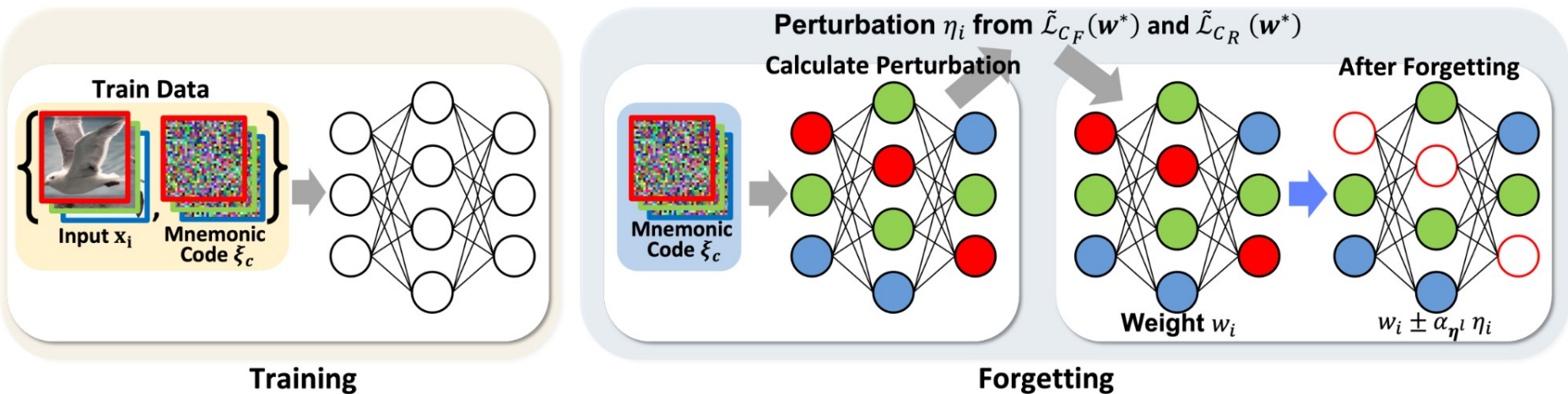
$$H(\theta) = \begin{bmatrix} \frac{\partial^2 f}{\partial \theta_1^2} & \frac{\partial^2 f}{\partial \theta_1 \partial \theta_2} & \cdots \\ \frac{\partial^2 f}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 f}{\partial \theta_2^2} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

Diagonal Approximation

$$H(\theta) = \begin{bmatrix} \frac{\partial^2 f}{\partial \theta_1^2} & \frac{\partial^2 f}{\partial \theta_1 \partial \theta_2} & \cdots \\ \frac{\partial^2 f}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 f}{\partial \theta_2^2} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

$$H_{\text{diag}}(\theta) = \begin{bmatrix} \frac{\partial^2 f}{\partial \theta_1^2} & 0 & \cdots \\ 0 & \frac{\partial^2 f}{\partial \theta_2^2} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

Method of the Paper (1/5)



$$\xi \sim N(0, 1)$$

Method of the Paper (2/5)

Algorithm 1 Training with mnemonic code

Input: dataset $\mathbf{x} \sim p^d(\mathbf{x})$, model parameter \mathbf{w} , loss \mathcal{L}

Parameter: mnemonic code replacing probability t_{mix} ,
learning rate lr

Output: trained model parameters

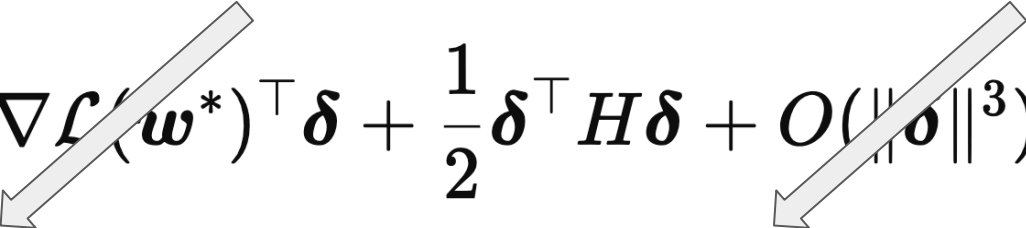
```
1:  $\xi \sim N(\mathbf{0}, \mathbf{1})$ 
2: for e in epochs do
3:   for i in datasize do
4:      $t \sim U(0, 1)$ 
5:     if  $t < t_{\text{mix}}$  then
6:        $\tilde{\mathbf{x}}_i = \xi_c$ 
7:     else
8:        $\tilde{\mathbf{x}}_i = \mathbf{x}_i$ 
9:     end if
10:  end for
11:   $\mathbf{w} = \mathbf{w} - \text{lr} \nabla_{\mathbf{w}} \mathcal{L}(\tilde{\mathbf{x}}; \mathbf{w})$ 
12: end for
```

Method of the Paper (3/5)

$$\begin{aligned} & \mathcal{L}_{\mathcal{C}_F}(\boldsymbol{w}^* + \boldsymbol{\delta}) \\ & \simeq \mathcal{L}_{\mathcal{C}_F}(\boldsymbol{w}^*) + \frac{1}{2} \boldsymbol{\delta}^T F_{\mathcal{C}_F} \boldsymbol{\delta} \\ & = \mathcal{L}_{\mathcal{C}_F}(\boldsymbol{w}^*) + \frac{1}{2} \boldsymbol{\delta}^T \{t_{\text{mix}} F_{\mathcal{C}_F}^{\xi} + (1 - t_{\text{mix}}) F_{\mathcal{C}_F}^d\} \boldsymbol{\delta} \\ & \simeq \mathcal{L}_{\mathcal{C}_F}(\boldsymbol{w}^*) + \frac{1}{2} \{t_{\text{mix}} \sum_i \boldsymbol{f}_{\mathcal{C}_F, i}^{\xi} + (1 - t_{\text{mix}}) \sum_i \boldsymbol{f}_{\mathcal{C}_F, i}^d\} \boldsymbol{\delta}_i^2, \end{aligned} \tag{2}$$

Taylor Series Expansion / Laplace's Method

$$f(x) = f(a) + f'(a)(x - a) + \frac{1}{2}f''(a)(x - a)^2 + O((x - a)^3)$$

$$\mathcal{L}(\mathbf{w}^* + \boldsymbol{\delta}) \approx \mathcal{L}(\mathbf{w}^*) + \nabla \mathcal{L}(\mathbf{w}^*)^\top \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^\top H \boldsymbol{\delta} + O(\|\boldsymbol{\delta}\|^3)$$


$$\mathcal{L}(\mathbf{w}^* + \boldsymbol{\delta}) \approx \mathcal{L}(\mathbf{w}^*) + \frac{1}{2} \boldsymbol{\delta}^\top H \boldsymbol{\delta}$$

$$\mathcal{L}_{\mathcal{C}_F}(\mathbf{w}^* + \boldsymbol{\delta}) \approx \mathcal{L}_{\mathcal{C}_F}(\mathbf{w}^*) + \frac{1}{2} \sum_i f_{\mathcal{C}_F, i} \delta_i^2$$

Method of the Paper (4/5)

$$w_i = w_i \pm \alpha_{\boldsymbol{\eta}^l} \eta_i, \quad (5)$$

$$\eta_i = \frac{f_{\mathcal{C}_F, i}}{f_{\mathcal{C}_R, i}} = \frac{\frac{1}{|\#\mathcal{C}_F|} \sum_{j \in \mathcal{C}_F} \mathbb{E} \left[\left(\frac{\partial \mathcal{L}_j}{\partial w_i} \right)^2 \right]}{\frac{1}{|\#\mathcal{C}_R|} \sum_{k \in \mathcal{C}_R} \mathbb{E} \left[\left(\frac{\partial \mathcal{L}_k}{\partial w_i} \right)^2 \right]}, \quad (6)$$

$$\alpha_{\boldsymbol{\eta}^l} = \min \left(\lambda_1, \frac{\lambda_2}{\max_{\eta_i \in \boldsymbol{\eta}^l} \eta_i} \right), \quad (7)$$

Metho

Algorithm 2 Forgetting with mnemonic code

Input: trained model parameter \mathbf{w} , loss \mathcal{L} , forget class set \mathcal{C}_F , remain class set \mathcal{C}_R , mnemonic codes ξ , layers $\{l_1, l_2, \dots\}$

Parameter: λ_1, λ_2

Output: Forgotten parameters

```
1:  $\mathbf{f}_{\mathcal{C}_F} = \mathbf{0}$ 
2:  $\mathbf{f}_{\mathcal{C}_R} = \mathbf{0}$ 
3: for  $c$  in  $\mathcal{C}_F$  do
4:    $\mathbf{f}_{\mathcal{C}_F} = \mathbf{f}_{\mathcal{C}_F} + \nabla_{\mathbf{w}} \mathcal{L}(\xi_c; \mathbf{w})$ 
5: end for
6: for  $c$  in  $\mathcal{C}_R$  do
7:    $\mathbf{f}_{\mathcal{C}_R} = \mathbf{f}_{\mathcal{C}_R} + \nabla_{\mathbf{w}} \mathcal{L}(\xi_c; \mathbf{w})$ 
8: end for
9:  $\mathbf{f}_{\mathcal{C}_F} = \mathbf{f}_{\mathcal{C}_F} / |\mathcal{C}_F|$ 
10:  $\mathbf{f}_{\mathcal{C}_R} = \mathbf{f}_{\mathcal{C}_R} / |\mathcal{C}_R|$ 
11:  $\boldsymbol{\eta} = \frac{\mathbf{f}_{\mathcal{C}_F}}{\mathbf{f}_{\mathcal{C}_R}}$ 
12: for  $l$  in layers do
13:    $\alpha_{\boldsymbol{\eta}^l} = \min \left( \lambda_1, \frac{\lambda_2}{\max_{\boldsymbol{\eta}^l \in \boldsymbol{\eta}^l} \boldsymbol{\eta}^l} \right)$ 
14:    $\mathbf{w}_1^l = \mathbf{w}^l + \alpha_{\boldsymbol{\eta}^l} \boldsymbol{\eta}^l$ 
15:    $\mathbf{w}_2^l = \mathbf{w}^l - \alpha_{\boldsymbol{\eta}^l} \boldsymbol{\eta}^l$ 
16: end for
17: if  $A_R(\mathbf{w}_1) + E_F(\mathbf{w}_1) > A_R(\mathbf{w}_2) + E_F(\mathbf{w}_2)$  then
18:   return  $\mathbf{w}_1$ 
19: else
20:   return  $\mathbf{w}_2$ 
21: end if
```

Results (1/4)

Method	Processing Time	Data-Free	MU Target
CertifiedRemoval (Guo et al., 2020)	$\mathcal{O}(N_R + N_F)$	✗	item
SISA (Bourtole et al., 2021)	$\mathcal{O}(E \cdot \frac{N_R}{M})$	✗	item
Arcane (Yan et al., 2022)	$\mathcal{O}(E \cdot \frac{N_R}{C_R + C_F})$	✗	item
FastMU (Tarun et al., 2023)	$\mathcal{O}(S \cdot C_F + E \cdot C_F + N_R)$	✓	class
ZeroShotMU (Chundawat et al., 2023)	$\mathcal{O}((S + E)(C_F + C_R))$	✓	class
LwSF (Shibata et al., 2021)	$\mathcal{O}(E(N_{\text{new}} + C_R))$	✓	class/task
SFDN (Golatkhar et al., 2020a)	$\mathcal{O}(N_R)$	✗	class/item
NTK-F (Golatkhar et al., 2020b)	$\mathcal{O}(N_R + N_F)$	✗	class/item
SSD (Foster et al., 2024)	$\mathcal{O}(N_R + N_F)$	✗	class/item
ERM-KTP (Lin et al., 2023)	$\mathcal{O}(E \cdot N_R)$	✗	class
Ours	$\mathcal{O}(C_R + C_F)$	✓	class

Results (2/4)

Table 2: **Comparison results in A_R .** We evaluate the baseline and our methods three times and provide the mean and standard deviation. The highest values are shown in bold.

	MNIST	CIFAR10	CUB	STN
FastMU	96.5 \pm 0.1	90.4 \pm 0.5	73.1 \pm 1.3	88.0 \pm 0.1
LwSF	43.7 \pm 9.6	65.4 \pm 16.6	68.2 \pm 3.5	80.1 \pm 6.7
SFDN	94.1 \pm 0.7	93.4 \pm 0.2	78.2 \pm 0.6	88.3 \pm 0.6
SSD	96.9 \pm 0.0	94.2 \pm 0.0	44.3 \pm 0.0	74.4 \pm 0.0
ERM-KTP	-	92.7 \pm 0.4	42.8 \pm 3.2	75.6 \pm 4.0
Ours	95.9 \pm 0.1	94.4 \pm 0.1	79.3 \pm 0.7	91.7 \pm 0.3

Results (3/4)

Table 3: **Comparison results in E_F .** We evaluate the baseline and our methods three times and provide the mean and standard deviation. The highest values are shown in bold.

	MNIST	CIFAR10	CUB	STN
FastMU	98.0 \pm 0.3	100 \pm 0.0	68.6 \pm 12.0	60.9 \pm 6.9
LwSF	94.4 \pm 1.7	100 \pm 0.0	93.1 \pm 7.0	98.2 \pm 1.8
SFDN	100 \pm 0.0	96.3 \pm 2.5	100 \pm 0.0	100 \pm 0.0
SSD	93.1 \pm 0.0	100 \pm 0.0	100 \pm 0.0	100 \pm 0.0
ERM-KTP	-	100 \pm 0.0	100 \pm 0.0	100 \pm 0.0
Ours	100 \pm 0.0	100 \pm 0.0	100 \pm 0.0	100 \pm 0.0

Results (4/4)

Table 4: **Forgetting for ImageNet dataset.** We perform fine-tuning using mnemonic code on the pre-trained model and forget with our method. We show the forgetting capability and processing time for the pre-trained, fine-tuned, and forgotten models, respectively.

Architecture		$A_R \uparrow$	$E_F \uparrow$	Time [s] \downarrow
ResNet-18	Pretrained	69.8	12.0	-
	Fine-tuned	67.5	12.0	882
	After MU	67.5	100	8.66
ResNeXt-50	Pretrained	77.4	6.0	-
	Fine-tuned	75.9	12.0	6923
	After MU	75.9	100	21.0
Swin-Transformer	Pretrained	80.9	4.0	-
	Fine-tuned	78.8	4.0	8488
	After MU	75.3	92.0	28.6

Hyperparameter Search

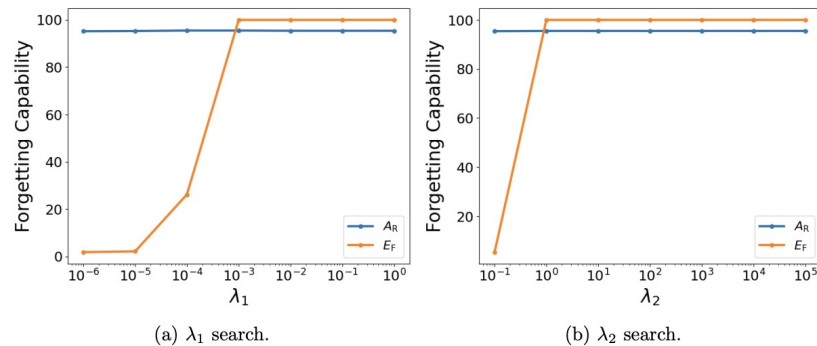


Figure 6: Hyperparameter search on MNIST.

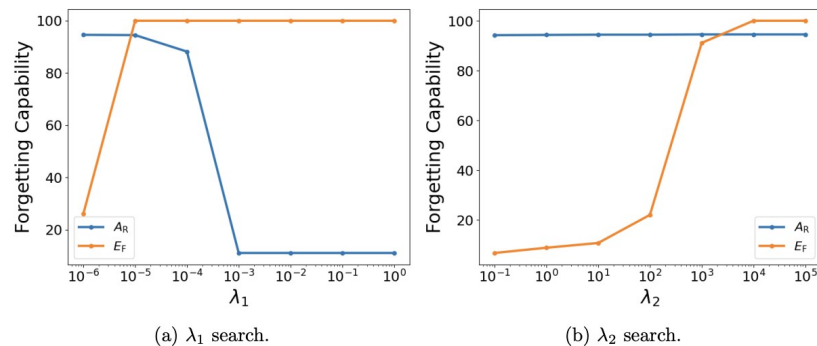


Figure 7: Hyperparameter search on CIFAR10.

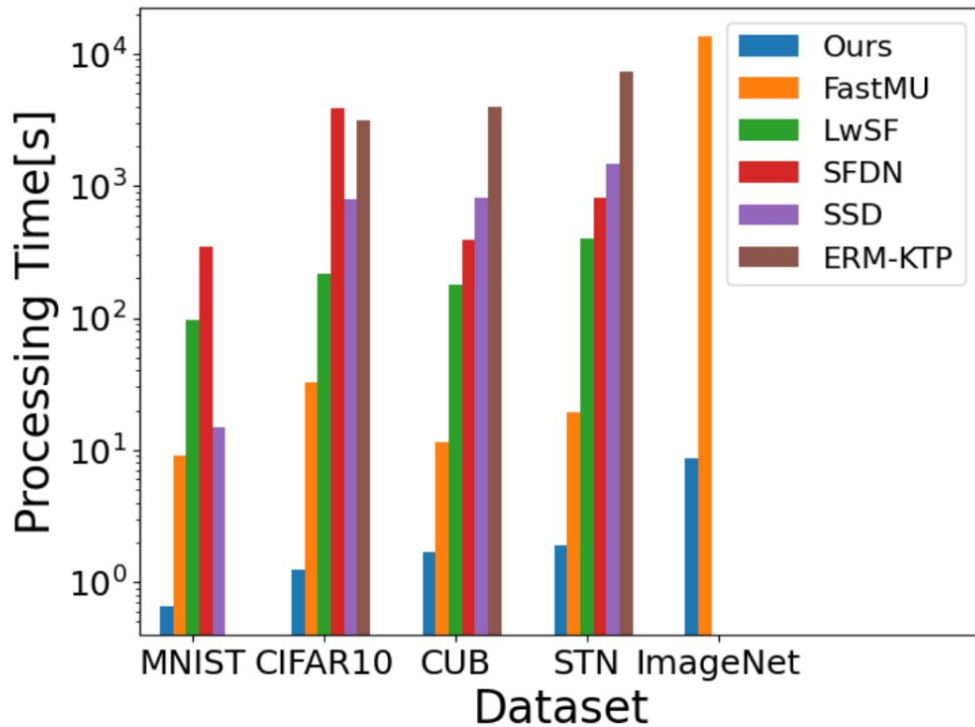


Figure 3: **Comparison results in MU processing time.** We measure the forgetting time concerning our method and the baselines.

Future Work

- Actually using Fisher Information Matrix
- Calculating Fisher Information only in important layers
- More realistic scenarios (not class-wise unlearning)