

## Appendix

In this Appendix, Appendix A provides experiment details and deferred experiment figures and tables; Appendix B provides the deferred proofs and the complete algorithm description. In addition, Appendix B also describes how to generalize *MUter* from successive *single datapoint* unlearning to successive *batch of datapoints* removal. Our source code can be found in Supplementary Material.

### A. Experiment Details and Deferred Experiment Figures and Tables

#### A.1. Experiment Details

Our experiments are based on the Pytorch platform and run on RTX 3090. For the pretraining the Wide ResNet model on Downsampled ImageNet, we run on four GPU devices; For all other experiments, we run on a single GPU device.

We summarize the hyperparameters used for our adversarial training/finetuning in Table 3. For Neural Network Model, we conduct experiments on Wide ResNet 28-10 model using the datasets Lacuna-10 and CIFAR-10. We first perform adversarial pretraining on Lacuna-100 and Downsampled ImageNet datasets, both by SGD with momentum for the outer-loops. For the former, we fix the learning rate to 0.1, and then train for 80 epochs with momentum 0.9 and weight decay 0.0005. For the latter, we use the pretrained model in [29] as the pretrained model here, which has a similar training process. Then, we freeze all but the last layer to adversarially finetune the model on Lacuna-10 and CIFAR-10 datasets.

We summarize the hyperparameters required by *MUter* in Table 4. In the neural network model experiments, [32] points out that under the non-convex setting, the Hessian matrix on parameters  $\omega^*$  sometimes will not be positive definite. We follow [32] to add a damping term  $\lambda$  on the diagonal with  $\lambda = 0.0001$ .

Dataset	Model type	Learning rate	(Tuning) Epoch	FGSM epsilon	PGD		
					epsilon	alpha	steps
MNIST-b	Logistic	0.01	100	0.25	0.25	0.03	15
	Ridge	0.01	100	0.25	0.25	0.03	15
Covtype	Logistic	0.1	15	4/255	4/255	0.004	7
	Ridge	0.1	15	4/255	4/255	0.004	7
Lacuna-10	Neural Network	0.01	20	8/255	8/255	2/255	10
CIFAR-10	Neural Network	0.001	10	8/255	8/255	2/255	10

Table 3: Adversarial training/finetuning parameters.

#### A.2. Datasets

**Linear Model.** We perform experiments on MNIST-b and Covtype. MNIST-b is taken from the MNIST dataset, which consists of  $28 * 28$  grayscale images from digit ‘0’ to digit ‘9’. We select the digit ‘1’ and digit ‘7’ to form the binary

Dataset	Model type	Neumann Series order $k$	Conjugate Gradient iterations $C$
MNIST-b	Logistic	3	10
	Ridge	3	10
Covtype	Logistic	3	20
	Ridge	20	20
Lacuna-10	Neural Network	100	10
CIFAR-10	Neural Network	100	20

Table 4: *MUter* parameters.

subset MNIST-b. Covtype with 54 attributes is used to classify the main tree species in the Roosevelt National Forest wilderness area, where use the binary classification version from LIBSVM.

**Neural Network Model.** We introduce Lacuna-100 and downsampled ImageNet as core datasets, and conduct experiments on target datasets Lacuna-10 and CIFAR-10. CIFAR-10 consists of  $32 * 32$  color pictures, covering different animals and machines in 10 categories. The Downsampled ImageNet is derived from the 1000-class ImageNet dataset, which resize the image to  $32 * 32$ . Lacuna-10/Lacuna-100 comes from [22]. We use the same data processing method to select 10/100 celebrities (no intersection) from VGGFace2 [7], and each celebrity randomly selects at least 500 pictures. Then each celebrity divides 100 pictures to form the test set, and the remaining images to form the training set. Finally, we resize the images to  $32 * 32$ .

To facilitate the constrained adversarial perturbation for adversarial training, all the above datasets are scaled to  $[0, 1]$  (for image data, we use Totensor to transform, and for Covtype, we choose the version scaled to  $[0, 1]$  in LIBSVM.). We summarize the dimensions, classes, and quantity information of the above datasets in Table 5.

Dataset	Domension	Classes	Train data	Test Data
MNIST-b	784	2	11,982	1,198
Covtype	54	2	522,910	58,102
Lacuna-10	3072	10	4,374	1,000
CIFAR-10	3072	10	50,000	10,000

Table 5: Datasets Statistics

#### A.3. Deferred Experiment Figures and Tables

In this section, we report additional experiment results, where Subsection A.3.1 reports additional experiment results under the PGD setting, and Subsection A.3.2 reports additional experiment results under the FGSM setting.

##### A.3.1 Deferred Experiment Figures under PGD Setting

We summarize the experiment results under Ridge Model with PGD in Figure 6.

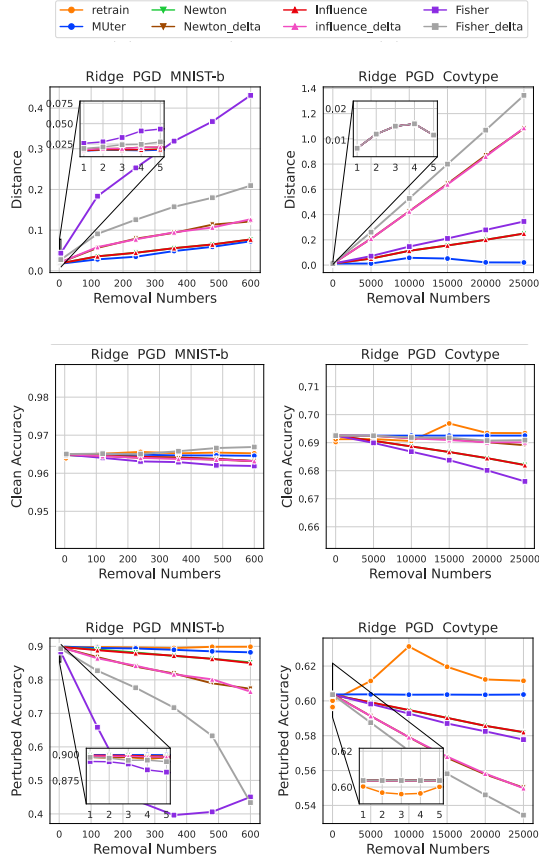


Figure 6: Evaluation results on Ridge Regression Model with PGD: Effectiveness (left column), Accuracy (middle column), and Robustness (right column) on datasets MNIST (top) and Covtype (bottom). Large plots have greater removal numbers: 1%, 2%, ..., 5%; Small plots inside the large plots have fewer removal numbers: 1, 2, ..., 5.

### A.3.2 Deferred Experiment Figures and Tables under FGSM Setting

In this part, we change the way of generating perturbations from PGD to FGSM. With the same experimental settings under PGD conditions, we report experiment results both on Linear Model and Neural Network Model.

**Results with Linear Model.** We report the effectiveness, accuracy and robustness metrics of logistic model and ridge model in Figure 7 and Figure 8, respectively. In Table 6, we summarize the efficiency comparison for linear model under FGSM setting.

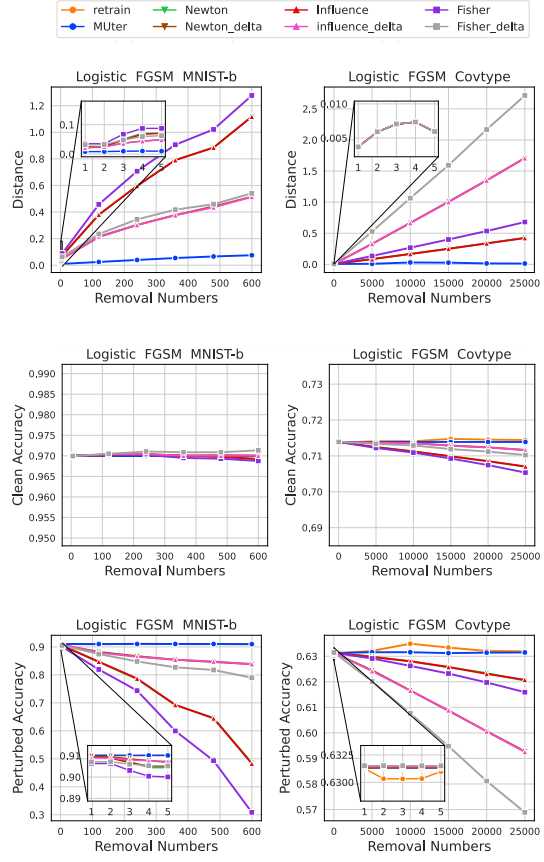


Figure 7: Evaluation results on Logistic Regression Model with FGSM: Effectiveness (left column), Accuracy (middle column), and Robustness (right column) on datasets MNIST (top) and Covtype (bottom). Large plots have greater removal numbers: 1%, 2%, ..., 5%; Small plots inside the large plots have fewer removal numbers: 1, 2, ..., 5.

Model Type	Removal Number	MNIST-b				Covtype			
		Fisher	F-delta	MUter	Retrain	Fisher	F-delta	MUter	Retrain
LR	1	0.002	0.002	0.008	14.6	0.002	0.002	0.007	75
	5	0.010	0.012	0.045	14.5	0.008	0.009	0.033	76
	10	0.020	0.022	0.089	14.2	0.018	0.019	0.065	77
RR	1%	0.242	0.249	0.978	14.3	9.178	11.278	33.974	77
	1	0.002	0.002	0.008	14.6	0.002	0.002	0.007	78
	5	0.009	0.012	0.047	14.5	0.009	0.009	0.034	78
	10	0.018	0.020	0.091	14.3	0.018	0.021	0.063	79
	1%	0.239	0.244	0.993	14.5	9.797	11.043	33.505	79

Table 6: Efficiency results with Logistic Regression Model (top) and Ridge Regression Model (bottom) under FGSM: The unlearning time (in seconds) of **Fisher**, **Fisher-delta**, **MUter** and **Retrain** under varying removal numbers: 1, 5, 10, 1% (120 for MNSIT-b, 5000 for Covtype).

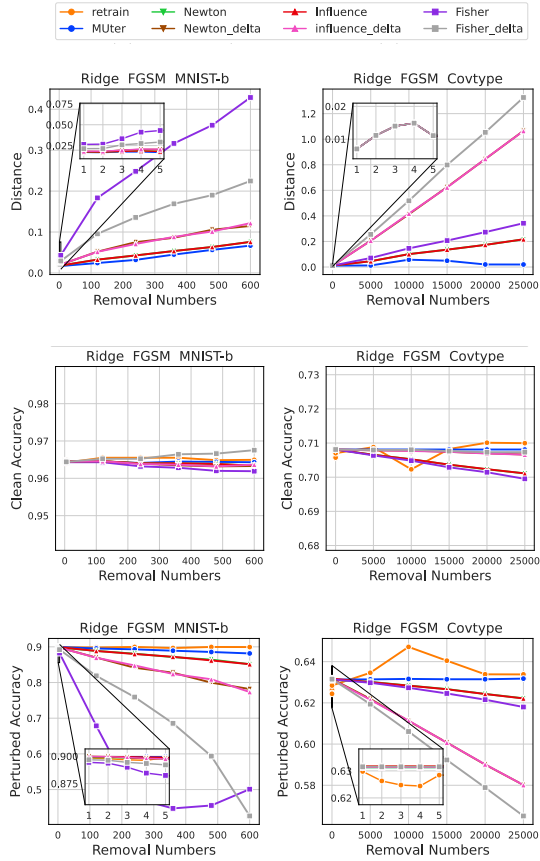


Figure 8: Evaluation results on Ridge Regression Model with FGSM: Effectiveness (left column), Accuracy (middle column), and Robustness (right column) on datasets MNIST (top) and Covtype (bottom). Large plots have greater removal numbers: 1%, 2%, ..., 5%; Small plots inside the large plots have fewer removal numbers: 1, 2, ..., 5.

**Result with Neural Network Model.** We report the effectiveness, accuracy, and robustness metrics of the neural network model in Figure 9. In Table 7, we summarize the efficiency comparison for the neural network model under the FGSM setting.

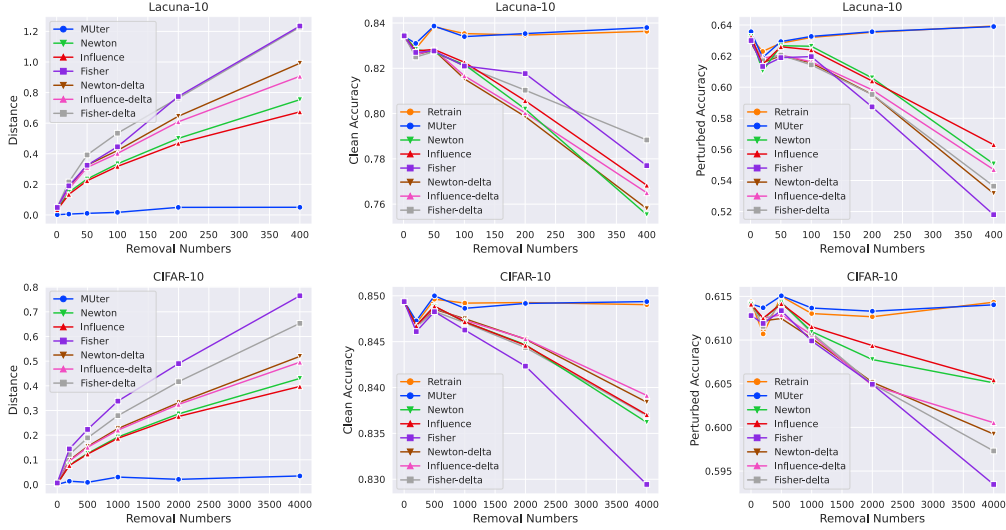


Figure 9: Evaluation results on Neural Network with FGSM: Effectiveness (left column), Accuracy (middle column), and Robustness (right column) on datasets Lacuna-10 (top) and CIFAR-10 (bottom), under removal numbers (1, 0.4%, 1%, 2%, 4%, 8%).

Removal Number	Lacuna-10				CIFAR-10			
	Fisher	F-delta	MUTEr	Retrain	Fisher	F-delta	MUTEr	Retrain
1	1.51	1.57	3.96	150	1.43	1.56	3.72	799
5	7.81	8.08	20.72	151	7.71	8.11	18.23	794
10	15.46	15.94	39.82	150	15.42	15.87	36.84	790

Table 7: Efficiency results with Neural Network Model under FGSM: The unlearning time (in seconds) of **Fisher**, **Fisher-delta**, **MUTEr** and **Retrain** under varying removal numbers: 1, 5, 10.

## B. Deferred Proofs and Algorithm Description

In this section, Subsection B.1 provides the omitted proofs for the derivation of the ATM unlearning update (i.e., Lemma 1 and Theorem 1), Subsection B.2 provides the omitted proof for the derivation of the successive unlearning setting (i.e., Corollary 1), Subsection B.3 provides the deferred Lemma for Schur complement for completeness and the omitted proof for Theorem 2, Subsection B.4 provides the complete algorithm description for *MUter*, and Subsection B.5 how to generalize *MUter* from successive *single datapoint* unlearning to successive *batch of datapoints* removal.

### B.1. Proofs for Derivation of ATM Unlearning Update

#### B.1.1 Proof of Lemma 1

*Proof.* By Taylor expansion around both  $\omega^*$  and  $\delta_i(\omega^*)$ , we have

$$\begin{aligned} & \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n \nabla_{\omega} l(\omega^u, \mathbf{x}_i + \delta_i(\omega^u)) \\ &= \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n \left[ \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) \right. \\ & \quad + \partial_{\omega} \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega^u - \omega^*) + O(\|\omega^u - \omega^*\|_2^2) \\ & \quad \left. + \partial_{\omega} \delta l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\delta_i(\omega^u) - \delta_i(\omega^*)) + O(\|\delta_i(\omega^u) - \delta_i(\omega^*)\|_2^2) \right]. \end{aligned} \quad (17)$$

By furthering neglecting the higher-order Taylor expansion terms  $O(\|\omega^u - \omega^*\|_2^2)$  and  $O(\|\delta_i(\omega^u) - \delta_i(\omega^*)\|_2^2)$ , we have the following relationship,

$$\begin{aligned} & \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n \nabla_{\omega} l(\omega^u, \mathbf{x}_i + \delta_i(\omega^u)) \\ & \approx \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n \left[ \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) \right. \\ & \quad + \partial_{\omega} \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega^u - \omega^*) \\ & \quad \left. + \partial_{\omega} \delta l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\delta_i(\omega^u) - \delta_i(\omega^*)) \right], \end{aligned} \quad (18)$$

which proves Lemma 1. □

#### B.1.2 Proof of Theorem 1

*Proof.* By Lemma 1, we begin with

$$\begin{aligned} & \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n \nabla_{\omega} l(\omega^u, \mathbf{x}_i + \delta_i(\omega^u)) \\ & \approx \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n \left[ \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) \right. \\ & \quad \left. + \partial_{\omega} \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega^u - \omega^*) + \partial_{\omega} \delta l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\delta_i(\omega^u) - \delta_i(\omega^*)) \right]. \end{aligned} \quad (19)$$

For the first line on the right hand side of eq.(19), we have

$$\begin{aligned}
& \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) \\
&= \frac{1}{n-1} \sum_{i=1}^n \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) - \frac{1}{n-1} \nabla_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) \\
&= 0 - \frac{1}{n-1} \nabla_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) \\
&= -\frac{1}{n-1} \nabla_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)),
\end{aligned} \tag{20}$$

where the second equality is obtained as follows. Since  $\omega^*$  is the optimum of the following ATM,

$$\omega^* = \operatorname{argmin}_{\omega} \max_{\delta_i \in \mathcal{B}(\mathbf{x}_i, r)} \frac{1}{n} \sum_{i=1}^n l(\omega, \mathbf{x}_i + \delta(\omega^*)), \tag{21}$$

it satisfies the following by Danskin's Theorem,

$$\frac{1}{n} \sum_{i=1}^n \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) = 0, \tag{22}$$

which is the second equality relation in eq.(20).

For the second line on the right hand side of eq.(19), we further expand  $(\delta_i(\omega^u) - \delta_i(\omega^*))$  by the implicit function theorem. That is, for all  $i \in \{1, \dots, n\}$ , we have

$$\delta_i(\omega^u) - \delta_i(\omega^*) = -\partial_{\delta}^{-1} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) \cdot \partial_{\delta} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) \cdot (\omega^u - \omega^*), \tag{23}$$

which gives

$$\begin{aligned}
& \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n \left[ \partial_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega^u - \omega^*) + \partial_{\delta} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\delta_i(\omega^u) - \delta_i(\omega^*)) \right] \\
&= \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n \left[ \partial_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega^u - \omega^*) \right. \\
&\quad \left. - \partial_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) \partial_{\delta}^{-1} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) \partial_{\delta} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega^u - \omega^*) \right] \\
&= \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n D_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega^u - \omega^*).
\end{aligned} \tag{24}$$

By substituting eq.(20) and eq.(24) into eq.(19), we have

$$\begin{aligned}
& \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n \nabla_{\omega} l(\omega^u, \mathbf{x}_i + \delta_i(\omega^u)) \\
&\approx -\frac{1}{n-1} \nabla_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) + \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n D_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega^u - \omega^*).
\end{aligned} \tag{25}$$

Thus,  $\omega^u = \omega^* + \mathbf{U}(\omega^*, i^\dagger) = \omega^* + \left[ \sum_{i=1, i \neq i^\dagger}^n D_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) \right]^{-1} \cdot \left[ \nabla_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) \right]$  is the the solution to the following linear system,

$$-\frac{1}{n-1} \nabla_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) + \frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n D_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega^u - \omega^*) = 0, \tag{26}$$

which gives  $\frac{1}{n-1} \sum_{i=1, i \neq i^\dagger}^n \nabla_{\omega} l(\omega^u, \mathbf{x}_i + \delta_i(\omega^u)) \approx 0$  according to eq.(25). As a result, we have proved that  $\omega^u$  satisfies the approximate unlearning criteria for ATM.  $\square$

## B.2. Proofs for derivation of Successive Unlearning Setting

### B.2.1 Proof of Corollary 1

*Proof.* Denote the data that have been deleted by  $\mathcal{U}_r := \{i_1^\dagger, i_2^\dagger, \dots, i_r^\dagger\}$ . Similar to the ATM unlearning standard in eq.(3), we have the retraining-from-scratch model parameter after the  $r + 1$ -th unlearning by  $\omega_{-\mathcal{U}_r \cup i^\dagger}^*$ , which has the following definition,

$$\omega_{-\mathcal{U}_r \cup i^\dagger}^* = \operatorname{argmin}_{\omega} \frac{1}{n - r - 1} \sum_{i=1, i \neq \mathcal{U}_r \cup i^\dagger}^n \max_{\delta_i \in \mathcal{B}(\mathbf{x}_i, r)} l(\omega, \mathbf{x}_i + \delta_i). \quad (27)$$

The unlearning criteria for ATM at the  $r + 1$ -th unlearning is as follows,

$$\sum_{i=1, i \neq \mathcal{U}_r \cup i^\dagger}^n \nabla_{\omega} l(\omega_{-\mathcal{U}_r \cup i^\dagger}^*, \mathbf{x}_i + \delta_i(\omega_{-\mathcal{U}_r \cup i^\dagger}^*)) \approx 0. \quad (28)$$

Our aim is to show that  $\omega_{r+1}^u$  approximately satisfies the above criteria,

$$\sum_{i=1, i \neq \mathcal{U}_r \cup i^\dagger}^n \nabla_{\omega} l(\omega_{r+1}^u, \mathbf{x}_i + \delta_i(\omega_{r+1}^u)) \approx 0, \quad (29)$$

which is also the approximate unlearning criteria for ATM in eq.(6).

The proof is similar to Theorem 1, as follows,

$$\begin{aligned} & \sum_{i=1, i \neq \mathcal{U}_r \cup i^\dagger}^n \nabla_{\omega} l(\omega_{r+1}^u, \mathbf{x}_i + \delta_i(\omega_{r+1}^u)) \\ = & \left[ \sum_{i=1}^n \nabla_{\omega} l(\omega_{r+1}^u, \mathbf{x}_i + \delta_i(\omega_{r+1}^u)) \right] - \left[ \sum_{i=i_1^\dagger}^{i_r^\dagger} \nabla_{\omega} l(\omega_{r+1}^u, \mathbf{x}_i + \delta_i(\omega_{r+1}^u)) \right] - \left[ \nabla_{\omega} l(\omega_{r+1}^u, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega_{r+1}^u)) \right] \\ \approx & \left[ \left( \sum_{i=1}^n \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) + \partial_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega_{r+1}^u - \omega^*) + \partial_{\omega} \delta l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\delta_i(\omega_{r+1}^u) - \delta_i(\omega^*)) \right) \right] \\ & - \left[ \left( \sum_{i=i_1^\dagger}^{i_r^\dagger} \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) + \partial_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega_{r+1}^u - \omega^*) + \partial_{\omega} \delta l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\delta_i(\omega_{r+1}^u) - \delta_i(\omega^*)) \right) \right] \\ & - \left[ \left( \nabla_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) + \partial_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) (\omega_{r+1}^u - \omega^*) + \partial_{\omega} \delta l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) (\delta_{i^\dagger}(\omega_{r+1}^u) - \delta_{i^\dagger}(\omega^*)) \right) \right] \\ = & \left[ \sum_{i=1}^n \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) + D_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega_{r+1}^u - \omega^*) \right] \\ & - \left[ \sum_{i=i_1^\dagger}^{i_r^\dagger} \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) + D_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega_{r+1}^u - \omega^*) \right] \\ & - \left[ \nabla_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) + D_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) (\omega_{r+1}^u - \omega^*) \right] \\ = & \left[ \sum_{i=1}^n D_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega_{r+1}^u - \omega^*) \right] \\ & - \left[ \sum_{i=i_1^\dagger}^{i_r^\dagger} \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) + D_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) (\omega_{r+1}^u - \omega^*) \right] \\ & - \left[ \nabla_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) + D_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) (\omega_{r+1}^u - \omega^*) \right], \end{aligned} \quad (30)$$

where the approximation equality is by the Taylor expansion and neglecting the higher-order terms, the second equality is by implicit function theorem and the definition of the total Hessian, and the third equality is by  $\sum_{i=1}^n \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) = 0$  by eq.(22).

Thus,  $\omega_{r+1}^u = \omega^* + \mathbf{u}_r(\omega^*, i^\dagger)$  is the the solution to the following linear system,

$$\begin{aligned} & - \sum_{i=i_1^\dagger}^{i_r^\dagger} \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) - \nabla_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) \\ & + \left[ \sum_{i=1}^n \mathbb{D}_{\omega\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) - \sum_{i=i_1^\dagger}^{i_r^\dagger} \mathbb{D}_{\omega\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*)) - \mathbb{D}_{\omega\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}(\omega^*)) \right] (\omega^u - \omega^*) = 0, \end{aligned} \quad (31)$$

which gives  $\sum_{i=1, i \neq i_r \cup i^\dagger}^n \nabla_{\omega} l(\omega_{r+1}^u, \mathbf{x}_i + \delta_i(\omega_{r+1}^u)) \approx 0$ . As a result, we have proved that  $\omega^u$  satisfies the approximate unlearning criteria for ATM in eq.(29) that is consistent with eq.(6).  $\square$

### B.3. Proof for Schur Complement Conversion

#### B.3.1 Additional Lemma of Schur Complement

**Lemma 2.** (Schur Complement Conversion) Let  $\mathbf{S} = \mathbf{H}_{11} - \mathbf{H}_{12}\mathbf{H}_{22}^{-1}\mathbf{H}_{21}$ . If  $\mathbf{H}_{22}$  and  $\mathbf{S}$  are invertible, then  $\mathbf{H}$  is invertible and the following relation holds,

$$\mathbf{H}^{-1} = \begin{bmatrix} \mathbf{H}_{11} & \mathbf{H}_{12} \\ \mathbf{H}_{21} & \mathbf{H}_{22} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{S}^{-1} & -\mathbf{S}^{-1}\mathbf{H}_{12}\mathbf{H}_{22}^{-1} \\ -\mathbf{H}_{22}^{-1}\mathbf{H}_{21}\mathbf{S}^{-1} & \mathbf{H}_{22}^{-1}\mathbf{H}_{21}\mathbf{S}^{-1}\mathbf{H}_{12}\mathbf{H}_{22}^{-1} \end{bmatrix}. \quad (32)$$

#### B.3.2 Proof of Theorem 2

*Proof.* Let

$$\mathbf{S} = \left[ \mathcal{M}_r[\mathbb{D}_{\omega\omega}] - \partial_{\omega\omega} l_{i^\dagger}^* \right] - \left[ -\partial_{\omega\delta} \partial_{\delta\delta}^{-1} \partial_{\delta\omega} l_{i^\dagger}^* \right], \quad (33)$$

which is the Schur complement of the following block matrix,

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}_{11} & \mathbf{H}_{12} \\ \mathbf{H}_{21} & \mathbf{H}_{22} \end{bmatrix} = \begin{bmatrix} \mathcal{M}_r[\mathbb{D}_{\omega\omega}] - \partial_{\omega\omega} l_{i^\dagger}^* & \partial_{\omega\delta} l_{i^\dagger}^* \\ \partial_{\delta\omega} l_{i^\dagger}^* & -\partial_{\delta\delta} l_{i^\dagger}^* \end{bmatrix}. \quad (34)$$

Then, based on Lemma 2, we have

$$\mathbf{S}^{-1} \mathbf{g} = \Delta\omega = \begin{bmatrix} \mathbf{I} & \mathbf{0} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{H}_{11} & \mathbf{H}_{12} \\ \mathbf{H}_{21} & \mathbf{H}_{22} \end{bmatrix}^{-1} \cdot \begin{bmatrix} \mathbf{g} \\ \mathbf{0} \end{bmatrix}, \quad (35)$$

which can be cast as the solution of the following linear system,

$$\begin{bmatrix} \mathbf{H}_{11} & \mathbf{H}_{12} \\ \mathbf{H}_{21} & \mathbf{H}_{22} \end{bmatrix} \cdot \begin{bmatrix} \Delta\omega \\ \Delta\alpha \end{bmatrix} = \begin{bmatrix} \mathbf{g} \\ \mathbf{0} \end{bmatrix}. \quad (36)$$

By substituting eq.(34) in, we have

$$\begin{bmatrix} \mathcal{M}_r[\mathbb{D}_{\omega\omega}] - \partial_{\omega\omega} l_{i^\dagger}^* & \partial_{\omega\delta} l_{i^\dagger}^* \\ \partial_{\delta\omega} l_{i^\dagger}^* & -\partial_{\delta\delta} l_{i^\dagger}^* \end{bmatrix} \begin{bmatrix} \Delta\omega \\ \Delta\alpha \end{bmatrix} = \begin{bmatrix} \mathcal{M}_r[\nabla_{\omega}] + \nabla_{\omega} l_{i^\dagger}^* \\ \mathbf{0} \end{bmatrix}, \quad (37)$$

which proves Theorem 2.  $\square$



## B.4. The Complete Algorithm Description

We present the complete algorithm description for *MUter* in Algorithm 1.

---

### Algorithm 1 The Complete Algorithm Description for *MUter*

---

**Input:** Training dataset  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ , loss  $l$ , adversarial perturbation constraint  $\mathcal{B}(\mathbf{x}_i, r)$ , adversarial training finishing status  $\{\omega^*, \delta_1(\omega^*), \dots, \delta_n(\omega^*)\}$ , a sequence of indices to be removed  $\mathcal{U} = \{i_1^\dagger, i_2^\dagger, \dots\}$

1: **Stage I. Pre-Unlearning:**

2: Compute memory  $\mathcal{M}_0[\mathbb{D}\omega\omega] = \sum_{i=1}^n \widetilde{\mathbb{D}\omega\omega} l(\omega^*, \mathbf{x}_i + \delta_i(\omega^*))$  by eq.(16);

3: Initialize  $\mathcal{M}_0[\nabla\omega] = \mathbf{0}$  and  $\mathcal{M}_0[\omega^*] = \omega^*$ ;

4: **for**  $i_{r+1}^\dagger = i^\dagger$  **do**

5:   **Stage II. Unlearning:**

6:   Re-compute the adversarial perturbation  $\delta_{i^\dagger}(\omega^*)$ ;

7:   Compute gradient  $\mathbf{g} = \nabla_{\omega} l_{i^\dagger}^*$ ;

8:   Apply conjugate gradient (with  $C$  iterations) to solve the least square problem of the linear system

$$\begin{bmatrix} \mathcal{M}_r[\mathbb{D}\omega\omega] - \partial_{\omega\omega} l_{i^\dagger}^* & \partial_{\omega} \delta_{i^\dagger}^* \\ \partial_{\delta\omega} l_{i^\dagger}^* & -\partial_{\delta\delta} l_{i^\dagger}^* \end{bmatrix} \begin{bmatrix} \Delta\omega \\ \Delta\alpha \end{bmatrix} = \begin{bmatrix} \mathcal{M}_r[\nabla\omega] + \nabla_{\omega} l_{i^\dagger}^* \\ \mathbf{0} \end{bmatrix},$$

9:   Obtain  $\mathbf{U}_r(\omega^*, i_{r+1}^\dagger) = \Delta\omega$ , the model parameter after the  $r + 1$ -th unlearning:  $\omega_{r+1}^u = \omega^* + \mathbf{U}_r(\omega^*, i_{r+1}^\dagger)$ ;

10:   **Stage III. Post-Unlearning:**

11:    $\mathbf{m}_{i^\dagger} = \partial_{\omega\omega} l_{i^\dagger}^* - \partial_{\omega} \delta_{i^\dagger}^* \left[ \sum_{j=0}^k (\mathbf{I} - \partial_{\delta\delta} l_{i^\dagger}^*)^j \right] \partial_{\delta\omega} l_{i^\dagger}^*$  by eq.(16);

12:   Update the memory  $\mathcal{M}_{r+1}[\mathbb{D}\omega\omega] = \mathcal{M}_r[\mathbb{D}\omega\omega] - \mathbf{m}_{i^\dagger}$  and  $\mathcal{M}_{r+1}[\nabla\omega] = \mathcal{M}_r[\nabla\omega] + \mathbf{g}$ ;

13: **end for**

---

## B.5. Extension to Successive Batch Unlearning

In this subsection, we show that our method can be generalized to the successive batch unlearning setting.

**Successive Batch Unlearning Setting.** Denote the index sets of datapoints that have already been forgotten at timestamp  $r$  by  $\mathcal{U}_1^\dagger, \dots, \mathcal{U}_r^\dagger$ , where each  $\mathcal{U}_r^\dagger$  contains a set of data indices. Let the set of datapoints to be forgotten at timestamp  $r + 1$  by  $\mathcal{U}_{r+1}^\dagger = \mathcal{U}^\dagger$ . Corollary 2 below extends Corollary 1 to support successive batch unlearning for ATM.

**Corollary 2.** *Considering the successive batch unlearning setting, let the machine unlearning update at the  $r + 1$ -th timestamp  $\mathbf{U}_r(\omega^*, \mathcal{U}^\dagger)$  take the following form*

$$\begin{aligned} \mathbf{U}_r(\omega^*, i_{r+1}^\dagger) := & \left\{ \overbrace{\left[ \sum_{i=1}^n \mathbb{D}\omega\omega l(\omega^*, \mathbf{x}_i + \delta_i^*) \right]}^{\mathcal{M}_0[\mathbb{D}\omega\omega]} \right. \\ & \left. - \overbrace{\left[ \sum_{j=1}^r \sum_{i \in \mathcal{U}_j^\dagger} \mathbb{D}\omega\omega l(\omega^*, \mathbf{x}_i + \delta_i^*) \right]}^{\text{Part of } \mathcal{M}_r[\mathbb{D}\omega\omega]} - \left[ \sum_{i^\dagger \in \mathcal{U}^\dagger} \mathbb{D}\omega\omega l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}^*) \right]} \right\}^{-1} \\ & \cdot \overbrace{\left[ \sum_{j=1}^r \sum_{i \in \mathcal{U}_j^\dagger} \nabla_{\omega} l(\omega^*, \mathbf{x}_i + \delta_i^*) \right]}^{\mathcal{M}_r[\nabla\omega]} + \sum_{i^\dagger \in \mathcal{U}^\dagger} \nabla_{\omega} l(\omega^*, \mathbf{x}_{i^\dagger} + \delta_{i^\dagger}^*) \end{aligned} \quad (38)$$

Then, the unlearning model with updated parameters  $\omega_{r+1}^u := \omega^* + \mathbf{U}_r(\omega^*, \mathcal{U}^\dagger)$  satisfies the approximate unlearning criteria for the adversarial training model in eq.(6).

The proof of Corollary 2 is similar to the proof of Corollary 1. Based on Corollary 2, *MUter* can be similarly designed to support the successive batch unlearning setting.