

# Recalling The Forgotten Class Memberships: Unlearned Models Can Be Noisy Labelers to Leak Privacy

Zhihao Sui<sup>1</sup>, Liang Hu<sup>2\*</sup>, Jian Cao<sup>1\*</sup>, Dora D. Liu<sup>2</sup>

Usman Naseem<sup>3</sup>, Zhongyuan Lai<sup>4</sup>, Qi Zhang<sup>2</sup>

<sup>1</sup>Shanghai Jiao Tong University

<sup>2</sup>Tongji University

<sup>3</sup>Macquarie University

<sup>4</sup>Shanghai Ballsnow Intelligent Technology Co. Ltd

fancyboy@sjtu.edu.cn, lianghu@tongji.edu.cn, cao-jian@sjtu.edu.cn, liudongmei\_0506@163.com  
usman.naseem@mq.edu.au, abrikosoff@yahoo.com, zhangqi\_cs@tongji.edu.cn

## Abstract

Machine Unlearning (MU) technology facilitates the removal of the influence of specific data instances from trained models on request. Despite rapid advancements in MU technology, its vulnerabilities are still underexplored, posing potential risks of privacy breaches through leaks of ostensibly unlearned information. Current limited research on MU attacks requires access to original models containing privacy data, which violates the critical privacy-preserving objective of MU. To address this gap, we initiate the innovative study on recalling the forgotten class memberships from unlearned models (ULMs) without requiring access to the original one. Specifically, we implement a Membership Recall Attack (MRA) framework with a teacher-student knowledge distillation architecture, where ULMs serve as noisy labelers to transfer knowledge to student models. Then, it is translated into a Learning with Noisy Labels (LNL) problem for inferring correct labels of the forgetting instances. Extensive experiments on state-of-the-art MU methods with multiple real datasets demonstrate that the proposed MRA strategy exhibits high efficacy in recovering class memberships of unlearned instances. As a result, our study and evaluation have established a benchmark for future research on MU vulnerabilities.

## 1 Introduction

Data privacy is a core concept in the era of big data and extensive interconnectivity [Wu *et al.*, 2024; Bao *et al.*, 2023]. If a machine learning model has been trained on sensitive and private data, it can lead to significant security risks. In this context, the emergence of machine unlearning (MU) has been driven by stringent data privacy regulations such as GDPR [Hoofnagle *et al.*, 2019] and CCPA [Itakura and Terada, 2018], which require the removal of specific sensitive data

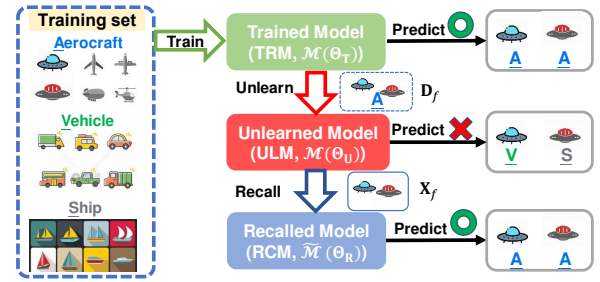


Figure 1: The demonstration of recall attack on the unlearned model (ULM) to recover class memberships having been unlearned via MU. This study is critical to the vulnerability of current MU models.

upon request. MU is designed to forget particular data points from the learned models [Cao and Yang, 2015]. As concerns about the increasing data misuse and privacy breaches, MU has gained more attention as a critical component in building safe machine learning systems.

Rapid advancements in MU research have increasingly posed a potential risk of privacy breaches by recovering the unlearned information about private data, highlighting the limited research on the full scope of MU vulnerabilities. In fact, the most relevant research [Hu *et al.*, 2024] that investigates inversion attacks against MU models was published recently. However, this work is based on an impractical assumption that unlearning inversion attacks require access to both **Trained Model (TRM,  $\mathcal{M}(\Theta_T)$ )**, and **Unlearned Model (ULM,  $\mathcal{M}(\Theta_U)$ )**, as shown in Figure 1. In general, only ULM is accessible to users, where sensitive privacy has been removed from TRM.

The membership inference attack (MIA) [Shokri *et al.*, 2017] is originally used to detect data samples used to train a machine learning model. Recently, MIA has been used to assess whether the influence of a forgetting dataset  $\mathcal{D}_f = \{\mathcal{X}_f, \mathcal{Y}_f\}$  has been successfully erased after MU [Chen *et al.*, 2021]. To go one step further, we formulate the attack on recovering the forgotten class memberships by only accessing ULM. As shown in Figure 1, our objective is to learn a **Recalled Model (RCM,  $\tilde{\mathcal{M}}(\Theta_R)$ )** from the ULM  $\mathcal{M}(\Theta_U)$  to correctly infer the labels of the input data  $\mathcal{X}_p$ , where  $\mathcal{X}_p$

\*Corresponding author

contains a subset of instances in  $\mathcal{X}_f$  in which the class memberships,  $\mathcal{Y}_f$ , have been forgotten in the ULM. As a result, we propose a Membership Recall Attack (MRA) framework to learn RCM in this paper. As most MU models restrict the scope of the investigation to the area of image classification tasks [Bourtoule *et al.*, 2021; Fan *et al.*, 2024; Jia *et al.*, 2023; Chen *et al.*, 2023; Chen *et al.*, 2024], we correspondingly study the proposed MRA on these MU models in this area.

To implement the MRA framework, we designed a teacher-student architecture to distill the knowledge from the ULM  $\mathcal{M}(\Theta_U)$  (as a teacher) to a **Student Model (STM,  $\tilde{\mathcal{M}}(\Theta_S)$ )**. More specifically, given an input image set  $\mathcal{X}_p$ , the ULM  $\mathcal{M}(\Theta_U)$  outputs the prediction labels  $\mathcal{M}(\mathcal{X}_p; \Theta_U) \mapsto \mathcal{Y}_p$ , where the prediction labels  $\mathcal{Y}_p$  may be noisy, especially when  $\mathcal{X}_p$  contains many instances of forgetting data  $\mathcal{X}_f$ . Therefore, we use the ULM  $\mathcal{M}(\Theta_U)$  to serve as a **noisy labeling teacher** for knowledge distillation, that is, using  $\{\mathcal{X}_p, \mathcal{Y}_p\}$  to train the STM  $\tilde{\mathcal{M}}(\Theta_S)$ . In this context, MRA can be further translated into a Learning with Noisy Labels (LNL) problem [Algan and Ulusoy, 2021] over  $\{\mathcal{X}_p, \mathcal{Y}_p\}$  that aims to infer correct labels from the noisy ones. Consequently, we selected samples with high confidence agreement between the teacher model  $\mathcal{M}(\Theta_U)$  and the STM  $\tilde{\mathcal{M}}(\Theta_S)$  for LNL. In particular, we discuss two cases for MRA, one is the closed-source case where the parameters of ULM  $\Theta_U$  are not accessible, and the other is the open-source case where the parameters are open for use. Moreover, we design a unified learning scheme of MRA to train the RCMs for these two cases. We summarize our contributions as follows.

- To our knowledge, this is the *first attempt* study of the recall attack of class membership, which can effectively assess the risk of data privacy breaches and promote the robustness of the MU study.
- We propose MRA, a model-agnostic attack framework, to effectively recover class memberships of forgotten instances from unlearned ULMs via various MU methods.
- We implement MRA with a teacher-student architecture where the ULM serves as a noisy labeling teacher to distill the knowledge to train the STM with noisy labels.
- We conducted extensive experiments in four widely used datasets in MU research, demonstrating both the theoretical and practical efficacy of our MRA approach against various SOTA MU methods.

## 2 Related Work

### 2.1 Machine Unlearning

**Exact Unlearning.** Retraining the model from scratch after removing specific data can intuitively and effectively achieve exact unlearning. In addition, [Bourtoule *et al.*, 2021] proposed SISA (Sharded, Isolated, Sliced, Aggregated) training, which trains isolated models on data shards for efficient unlearning by retraining only affected shards. Although effective, these unlearning approaches are computationally expensive and impractical for large-scale models and datasets.

**Approximate Unlearning.** The idea of modestly sacrificing the accuracy of forgetting in exchange for significant

improvements in unlearning efficiency has spurred the exploration of approximate unlearning techniques. Gradient ascent (**GA**) [Graves *et al.*, 2021; Golatkar *et al.*, 2020; Thudi *et al.*, 2022; Miao *et al.*, 2024] reverses the training of the model by adding gradients, thus moving the model towards greater loss for the data points targeted for removal. Random labeling (**RL**) [Golatkar *et al.*, 2020] that involves finetuning the original model on the forgetting dataset using random labels to enforce unlearning. Several methods estimate the impact of forgetting samples on the model parameters and conduct forgetting through the fisher information matrix (**FF**) [Becker and Liebig, 2022] or influence function (**IU**) [Koh and Liang, 2017; Izzo *et al.*, 2021].  $\ell_1$ -sparse (**L1-SP**) [Jia *et al.*, 2023] infuses weight sparsity into unlearning.

Moreover, most MU methods may degrade model performance, and lead to “*over-unlearning*”. Some recent work has explored more precise unlearning on target forget instances. Boundary unlearning (**BU**) [Chen *et al.*, 2023] shifts the decision boundary of the original model to imitate the decision behavior of the model retrained from scratch. **SalUn** [Fan *et al.*, 2024] introduces the concept of ‘weight saliency’ to narrow the performance gap with exact unlearning. To avoid over-unlearning, **UNSC** [Chen *et al.*, 2024] constrains the unlearning process within a null space tailored to the remaining samples to ensure that unlearning does not negatively impact the model performance.

### 2.2 Attacks on Machine Unlearning

Despite advances in MU techniques, the study of their vulnerabilities remains underexplored. To date, very limited MU attack methods have been proposed to affect efficiency [Marchant *et al.*, 2022] or fidelity [Di *et al.*, 2022; Hu *et al.*, 2023]. Studying attacks on MU is crucial to developing robust and secure MU methods.

**Membership Inference Attack (MIA).** MIA is originally used to infer if data samples are used to train a machine learning model [Shokri *et al.*, 2017]. With the development of MU, MIA has been widely used to check if the influence of forgetting data had been removed from the original model. However, [Chen *et al.*, 2021] show that MU can jeopardize privacy in terms of MIA. The goals of MIA and the proposed MRA are different. MIA aims to detect data samples if used for training, whereas MRA aims at recalling and inferring the class memberships of forgetting samples from ULMs.

**Model Inversion Attack.** It aims to reconstruct the original input data from the model outputs. [Fredrikson *et al.*, 2015] introduced model inversion attacks using the confidence scores output by a model to reconstruct input images. [Hu *et al.*, 2024] proposed the first inversion attack against unlearning. It extracts features and labels of forgetting samples, which most closely match the objectives of our study. Although the attack demonstrates notable effectiveness, it requires *access to the original TRM* before unlearning, which is impractical in real scenarios. In contrast, the proposed MRA only needs to access ULMs and supports more versatile MU methods. To our knowledge, we are the first to explore the attack **only using ULMs** to recall the class memberships of forgetting samples, without comparable prior work.

### 3 Preliminaries

We first introduce the datasets and models used in our study, followed by a formal definition of the problem.

#### 3.1 Involved Datasets

**Training dataset**  $\mathcal{D}_{tr} : \{\mathcal{X}_{tr}, \mathcal{Y}_{tr}\}$  is all data used to initially train machine learning models, where  $\mathcal{X}_{tr}$  denotes the image set and  $\mathcal{Y}_{tr}$  denotes the corresponding label set.

**Forgetting dataset**  $\mathcal{D}_f : \{\mathcal{X}_f, \mathcal{Y}_f\}$  is a subset of  $\mathcal{D}_{tr}$ , that is,  $\mathcal{D}_f \subset \mathcal{D}_{tr}$ . In MU,  $\mathcal{D}_f$  is a set of sensitive data that should be unlearned from the trained model, that is, the ULM cannot tell the true labels when  $\mathcal{X}_f$  is input.

**Remaining dataset**  $\mathcal{D}_r : \{\mathcal{X}_r, \mathcal{Y}_r\}$  is the remaining data of  $\mathcal{D}_{tr}$ , that is,  $\mathcal{D}_r = \mathcal{D}_{tr} \setminus \mathcal{D}_f$ , which should not be forgotten.

**Prediction dataset**  $\mathcal{D}_p : \{\mathcal{X}_p, \mathcal{Y}_p\}$  is the dataset for prediction, where  $\mathcal{X}_p = \mathcal{X}_{ts} \cup \mathcal{X}_u$  can be decomposed into two parts.  $\mathcal{X}_u \subseteq \mathcal{X}_f$  is the subset of the forgotten instances while  $\mathcal{X}_{ts}$  is the unseen dataset for testing.

#### 3.2 Involved Models

**Trained Model (TRM)**  $\mathcal{M}(\Theta_T)$  is the model that has been trained on the training dataset  $\mathcal{D}_{tr}$ .

**Unlearned Model (ULM)**  $\mathcal{M}(\Theta_U)$  is the model that has unlearned the forgetting dataset  $\mathcal{D}_f$  based on  $\mathcal{M}(\Theta_U)$ . It will serve as a noisy labeling teacher to distill knowledge.

**Student Model (STM)**  $\tilde{\mathcal{M}}(\Theta_S)$  is the model that receives the knowledge distilled from  $\mathcal{M}(\Theta_U)$ .

**Recalled Model (RCM)**  $\tilde{\mathcal{M}}(\Theta_R)$  is the model that has recalled forgotten class memberships based on  $\mathcal{M}(\Theta_U)$ .

#### 3.3 Problem Formulation

MU models remove the influence of forgetting the dataset  $\mathcal{D}_f$  from TRM  $\mathcal{M}(\Theta_T)$ , and release a ULM  $\mathcal{M}(\Theta_U)$  for public use. This paper aims to implement the MRA framework to recall forgotten class memberships given  $\mathcal{M}(\Theta_U)$ . In particular, we employ  $\mathcal{M}(\Theta_U)$  as a noisy labeler (i.e., teacher model) to distill the knowledge inferred from the prediction dataset  $\mathcal{X}_p$  to the STM  $\tilde{\mathcal{M}}(\Theta_S)$ .

Moreover, we discuss two common cases that lead to the final RCM  $\tilde{\mathcal{M}}(\Theta_R)$ . In the first case, ULM  $\mathcal{M}(\Theta_U)$  is usable but not trainable, e.g.  $\mathcal{M}(\Theta_U)$  is only accessible as a black-box service (**closed-source case**), and the STM  $\tilde{\mathcal{M}}(\Theta_S)$  will finally serve as RCM  $\tilde{\mathcal{M}}(\Theta_R)$ . In the second case,  $\mathcal{M}(\Theta_U)$  is trainable, for example,  $\mathcal{M}(\Theta_U)$  is released with its ULM parameters  $\Theta_U$  (**open-source case**), and the ULM  $\mathcal{M}(\Theta_U)$  is recovered to serve as the RCM  $\tilde{\mathcal{M}}(\Theta_R)$ .

## 4 Proposed Method

### 4.1 Overview

Firstly, we can easily obtain ULMs by applying various MU methods on a TRM. Given a ULM, Figure 2 demonstrates the workflow to implement the proposed MRA framework, which consists of two alternative learning steps to obtain the RCM  $\tilde{\mathcal{M}}(\Theta_R)$ .

**(1) Denosing Knowledge Distillation:** The ULM  $\mathcal{M}(\Theta_U)$  serves as a noisy labeler on the prediction dataset  $\mathcal{D}_p$  with

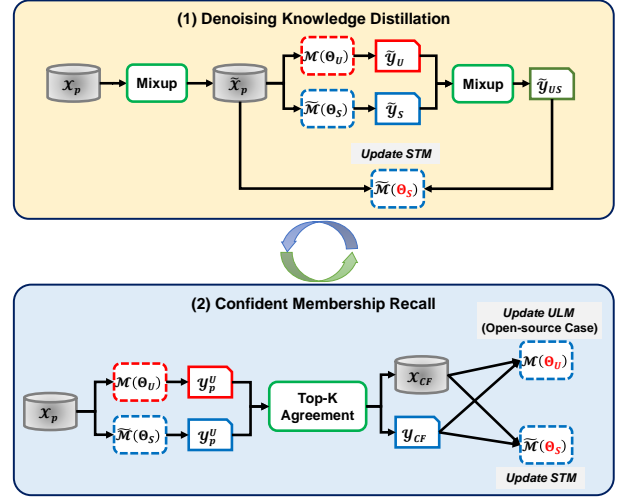


Figure 2: The workflow of proposed MRA framework. It consists of two alternative learning steps to obtain RCM: (1) Denosing Knowledge Distillation; (2) Confident Membership Recall.

some augmentation strategy, and the inferred pseudo labels with the augmented images are used to train STM  $\tilde{\mathcal{M}}(\Theta_S)$ .

**(2) Confident Membership Recall:** Both ULM  $\mathcal{M}(\Theta_U)$  and STM  $\tilde{\mathcal{M}}(\Theta_S)$  are to generate the prediction on  $\mathcal{D}_p$ , and the top- $K$  data samples with the highest probability of joint prediction, that is, the most confidently agreed pseudo labels for each class are selected to train STM  $\tilde{\mathcal{M}}(\Theta_S)$  and ULM  $\mathcal{M}(\Theta_U)$  (optional for the open source case).

### 4.2 Model Training and Unlearning

Given the training dataset  $\mathcal{D}_{tr} = \mathcal{D}_f \cup \mathcal{D}_r$ , we use  $\mathcal{D}_{tr}$  to train the model  $\mathcal{M}$ , which results in the TRM  $\mathcal{M}(\Theta_T)$  with the parameters  $\Theta_T$ . Then, we apply various SOTA MU methods on TRM  $\mathcal{M}(\Theta_T)$ , which leads to ULMs  $\mathcal{M}(\Theta_U)$ .

### 4.3 Implementation of MRA Framework

We implement the MRA framework with the following two alternative learning steps to obtain RCM  $\tilde{\mathcal{M}}(\Theta_R)$ .

#### (1) Denosing Knowledge Distillation

A sophisticated MU method should only unlearn the influence of forgetting dataset  $\mathcal{D}_f$  but retain the classification capability that is learned from the remaining dataset  $\mathcal{D}_r$ . As a result, the ULM  $\mathcal{M}(\Theta_U)$  can serve as a noisy labeler to distill knowledge to the STM  $\tilde{\mathcal{M}}(\Theta_S)$  given the set of prediction images  $\mathcal{X}_p$ .  $\mathcal{M}(\Theta_U)$  is prone to mislabeling on  $\mathcal{X}_p$  if it contains forgotten instances in  $\mathcal{D}_f$ . To avoid the input of original images that have been forgotten, we create augmented images by mixup [Zhang *et al.*, 2018] which is an effective regularization technique to deal with label noise [Carratino *et al.*, 2022]. Given an image  $x_1 \in \mathcal{X}_p$ , we randomly sample another image  $x_2$  from  $\mathcal{X}_p$ , and mix them as follows:

$$\tilde{x} = \beta_x \cdot x_1 + (1 - \beta_x) \cdot x_2 \quad (1)$$

where  $\beta_x \sim \text{Beta}(\alpha_x, \alpha_x)$  ( $\alpha_x = 0.2$  in this paper). Then, we take  $\tilde{x}$  as input to retrieve the soft pseudo label from ULM

$\mathcal{M}(\Theta_U)$  and STM  $\tilde{\mathcal{M}}(\Theta_S)$ :

$$\tilde{\mathbf{y}}_U = \mathcal{M}(\tilde{x}; \Theta_U), \quad \tilde{\mathbf{y}}_S = \tilde{\mathcal{M}}(\tilde{x}; \Theta_S) \quad (2)$$

Then, we can obtain the mixed soft pseudo label:

$$\tilde{\mathbf{y}}_{US} = \beta_y \cdot \tilde{\mathbf{y}}_U + (1 - \beta_y) \cdot \tilde{\mathbf{y}}_S \quad (3)$$

where  $\beta_y = 1$  is applied in the first warmup epoch (i.e. completely accept the knowledge from the teacher model) and  $\beta_y \sim \text{Beta}(\alpha_y, \alpha_y)$  ( $\alpha_y = 0.75$  in this paper) for label denoising after the warmup stage. Following Eqs (1) to (3), we can obtain  $\tilde{\mathcal{Y}}_S = \{\tilde{\mathbf{y}}_S\}$  and  $\tilde{\mathcal{Y}}_{US} = \{\tilde{\mathbf{y}}_{US}\}$  over the augmented set  $\tilde{x} \in \tilde{\mathcal{X}}_p$ . As a result, the parameters of STM  $\tilde{\mathcal{M}}(\Theta_S)$  can be updated by decreasing the mini-batch gradient in terms of minimizing cross-entropy (CE) loss.

$$\Theta_S = \arg \min_{\Theta_S} \text{CE}(\tilde{\mathcal{Y}}_S, \tilde{\mathcal{Y}}_{US}) \quad (4)$$

## (2) Confident Membership Recall

After the above step, STM  $\tilde{\mathcal{M}}(\Theta_S)$  has been trained on the augmented dataset based on  $\mathcal{X}_p$  with pseudo labels from the noisy labeler  $\mathcal{M}(\Theta_U)$ . Inspired by the LNL methods [Algan and Ulusoy, 2021], we design a balanced class membership recall strategy based on the highest confidence agreements between  $\mathcal{M}(\Theta_U)$  and  $\tilde{\mathcal{M}}(\Theta_S)$ . More specifically, we first input  $\mathcal{X}_p$  into both ULM  $\mathcal{M}(\Theta_U)$  and STM  $\tilde{\mathcal{M}}(\Theta_S)$  to predict soft labels (that is, probability over each class):

$$\mathbf{Y}_p^U = \mathcal{M}(\mathcal{X}_p; \Theta_U), \quad \mathbf{Y}_p^S = \tilde{\mathcal{M}}(\mathcal{X}_p; \Theta_S) \quad (5)$$

where  $\mathbf{Y}_p^U, \mathbf{Y}_p^S \in \mathbb{R}^{N \times C}$ ,  $N = |\mathcal{X}_p|$  denotes the number of samples in  $\mathcal{X}_p$  and  $C$  denotes the number of classes. For each instance  $x_i \in \mathcal{X}_p$ , we apply Laplace smoothing on its soft label  $\mathbf{y}_i^U = \mathbf{Y}_p^U[i, :]$  to avoid zero probability:

$$\tilde{\mathbf{y}}_i^U = \frac{\mathbf{y}_i^U + \gamma_l \cdot \mathbf{1}}{1 + C \cdot \gamma_l} \quad (6)$$

As a result, we obtain the smoothed probability matrices  $\tilde{\mathbf{Y}}_p^U$  over  $\mathcal{X}_p$ .  $\tilde{\mathbf{Y}}_p^S$  can be obtained in the same way. Then, we have the joint probability  $\tilde{\mathbf{y}}_i$  on  $x_i$ :

$$\tilde{\mathbf{y}}_i = \tilde{\mathbf{y}}_i^U \odot \tilde{\mathbf{y}}_i^S \quad \text{for } \tilde{\mathbf{y}}_i^U \in \tilde{\mathbf{Y}}_p^U, \tilde{\mathbf{y}}_i^S \in \tilde{\mathbf{Y}}_p^S \quad (7)$$

For all  $x_i \in \mathcal{X}_p$ , we have  $\tilde{\mathbf{Y}}_p = \tilde{\mathbf{Y}}_p^U \odot \tilde{\mathbf{Y}}_p^S$  in the matrix form, where  $\odot$  is the element-wise product. Given a class  $c$ , we have the joint probabilities  $\mathbf{y}_c = \tilde{\mathbf{Y}}_p[:, c]$  of all instances. A higher joint probability  $y_i \in \mathbf{y}_c$  implies greater confidence in the teacher and student models that the instance  $x_i$  should have the membership of the class  $c$ . Consequently, we apply a balanced strategy to select top- $K$  instances with the maximum joint probability for each class.

$$\mathcal{D}_{CF} = \{\mathcal{X}_{CF}, \mathcal{Y}_{CF}\} = \{(x_i, \tilde{\mathbf{y}}_c) | y_i \in \text{top-}K(\mathbf{y}_c) \quad (8)$$

$$\text{for } c \in \{1, \dots, C\}$$

where  $K = \lceil \tau \cdot N / C \rceil$ , and  $\tilde{\mathbf{y}}_c = (1 - \gamma_s) \cdot \mathbf{y}_c + \frac{\gamma_s}{K} \cdot \mathbf{1}$  denotes the smoothing of the label versus the hard label  $c$  ( $\mathbf{y}_c$  stands for the one-hot encoding of  $c$ ) which can effectively mitigate

## Algorithm 1 The Scheme of MRA

**Input:** Prediction Set:  $\mathcal{X}_p$ , Number of Epochs:  $K_E, K_D, K_R$

**Output:** Recalled Labels:  $\hat{\mathcal{Y}}_p$

```

1: for e in  $\{1, \dots, K_E\}$  do
2:   (1) Denosing Knowledge Distillation
3:   for i in  $\{1, \dots, K_D\}$  do
4:      $\tilde{\mathcal{X}}_p, \tilde{\mathcal{Y}}_{US} \triangleright$  Mixup augmentation by Eqs (1) to (3)
5:      $\Theta_S = \arg \min_{\Theta_S} \text{CE}(\tilde{\mathcal{Y}}_S, \tilde{\mathcal{Y}}_{US}) \triangleright$  Update STM
6:   end for
7:   (2) Confident Membership Recall
8:   for i in  $\{1, \dots, K_R\}$  do
9:      $\mathcal{D}_{CF} \triangleright$  Confident agreements by Eqs (5) to (8)
10:     $\Theta_S = \arg \min_{\Theta_S} \text{CE}(\mathcal{Y}_S, \mathcal{Y}_{CF}) \triangleright$  Update STM
11:    if  $\mathcal{M}(\Theta_U)$  is trainable (open-source case) then
12:       $\mathcal{D}_{CF} \triangleright$  Confident agreements by Eqs (5) to (8)
13:       $\Theta_U = \arg \min_{\Theta_U} \text{CE}(\mathcal{Y}_U, \mathcal{Y}_{CF}) \triangleright$  Update ULM
14:    end if
15:  end for
16: end for
17: return  $\hat{\mathcal{Y}}_p = \tilde{\mathcal{M}}(\mathcal{X}_p; \Theta_S) \triangleright$  Closed-source case
18:    $\hat{\mathcal{Y}}_p = \mathcal{M}(\mathcal{X}_p; \Theta_U) \triangleright$  Open-source case

```

label noise [Lukasik *et al.*, 2020]. Then, the STM  $\tilde{\mathcal{M}}(\Theta_S)$  is updated on  $\mathcal{D}_{CF}$  to learn the confident memberships.

$$\mathcal{Y}_S = \tilde{\mathcal{M}}(\mathcal{X}_{CF}; \Theta_S) \quad (9)$$

$$\Theta_S = \arg \min_{\Theta_S} \text{CE}(\mathcal{Y}_S, \mathcal{Y}_{CF}) \quad (10)$$

In the **open-source case**, the parameters of ULM  $\mathcal{M}(\Theta_U)$  are also accessible, so we will construct  $\mathcal{D}_{CF}$  to refine  $\Theta_U$  by Eqs (5) to (8) with the above updated STM  $\tilde{\mathcal{M}}(\Theta_S)$ .

$$\mathcal{Y}_U = \mathcal{M}(\mathcal{X}_{CF}; \Theta_U) \quad (11)$$

$$\Theta_U = \arg \min_{\Theta_U} \text{CE}(\mathcal{Y}_U, \mathcal{Y}_{CF}) \quad (12)$$

As a result, it leads to an alternative improvement co-training process between the teacher  $\mathcal{M}(\Theta_U)$  and the student  $\tilde{\mathcal{M}}(\Theta_S)$ , which can more effectively recall class memberships thanks to the knowledge retained by  $\mathcal{M}(\Theta_U)$ .

In Algorithm 1, we concisely summarize the above MRA scheme. After MRA, the updated STM  $\tilde{\mathcal{M}}(\Theta_S)$  (in closed-source case) or ULM  $\mathcal{M}(\Theta_U)$  (in open-source case) will serve as the RCM to predict the labels  $\hat{\mathcal{Y}}_p$  of  $\mathcal{X}_p$ .

## 5 Experiments

### 5.1 Experiment Setup

#### Data Preparation

In our experiments, four real datasets are used, which cover both low- and high-resolution data, providing a comprehen-

Dataset	# Classes	$\mathcal{D}_{tr}$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$
CIFAR-10	10	50,000	10,000	2500×5
CIFAR-100	100	50,000	10,000	250×5
Pet-37	37	3,680	3,669	50×5
Flower-102	102	3,074	1,020	224

Table 1: Statistic summary of the datasets and their splits

		TRM		FF	RL	GA	IU	BU	L1-SP	SalUn	UNSC
CIFAR-10 T: EFN S: EFN	$\mathcal{D}_{ts}$	0.833	ULM	0.164	0.388	0.505	0.105	0.462	0.451	0.486	0.193
			RCM	0.252	0.479	0.528	0.115	0.651	0.456	0.697	0.568
			$\Delta Acc$	<b>0.088</b>	<b>0.091</b>	<b>0.023</b>	<b>0.010</b>	<b>0.188</b>	<b>0.005</b>	<b>0.211</b>	<b>0.375</b>
	$\mathcal{D}_f$	1.000	ULM	0.142	0.153	0.267	0.090	0.082	0.154	0.117	0.031
			RCM	0.195	0.328	0.358	0.096	0.507	0.231	0.615	0.513
			$\Delta Acc$	<b>0.053</b>	<b>0.175</b>	<b>0.091</b>	<b>0.005</b>	<b>0.424</b>	<b>0.077</b>	<b>0.498</b>	<b>0.482</b>
CIFAR-100 T: EFN S: EFN	$\mathcal{D}_{ts}$	0.643	ULM	0.159	0.593	0.531	0.140	0.529	0.537	0.410	0.275
			RCM	0.295	0.607	0.551	0.297	0.566	0.550	0.522	0.480
			$\Delta Acc$	<b>0.137</b>	<b>0.014</b>	<b>0.020</b>	<b>0.157</b>	<b>0.037</b>	<b>0.013</b>	<b>0.112</b>	<b>0.206</b>
	$\mathcal{D}_f$	1.000	ULM	0.094	0.086	0.201	0.201	0.254	0.199	0.227	0.209
			RCM	0.214	0.537	0.270	0.256	0.423	0.258	0.358	0.372
			$\Delta Acc$	<b>0.119</b>	<b>0.451</b>	<b>0.070</b>	<b>0.055</b>	<b>0.169</b>	<b>0.058</b>	<b>0.130</b>	<b>0.163</b>
Pet-37 T: ResNet S: ResNet	$\mathcal{D}_{ts}$	0.895	ULM	0.257	0.754	0.740	0.622	0.646	0.740	0.706	0.769
			RCM	0.486	0.816	0.757	0.682	0.785	0.758	0.778	0.799
			$\Delta Acc$	<b>0.229</b>	<b>0.062</b>	<b>0.018</b>	<b>0.060</b>	<b>0.140</b>	<b>0.018</b>	<b>0.072</b>	<b>0.030</b>
	$\mathcal{D}_f$	1.000	ULM	0.224	0.332	0.224	0.200	0.084	0.196	0.128	0.304
			RCM	0.388	0.884	0.520	0.448	0.856	0.440	0.712	0.660
			$\Delta Acc$	<b>0.164</b>	<b>0.552</b>	<b>0.296</b>	<b>0.248</b>	<b>0.772</b>	<b>0.244</b>	<b>0.584</b>	<b>0.356</b>
Flower-102 T: Swin-T S: ResNet	$\mathcal{D}_{ts}$	0.939	ULM	0.294	0.756	0.314	0.512	0.706	0.530	0.496	NA
			RCM	0.376	0.831	0.510	0.589	0.758	0.655	0.599	NA
			$\Delta Acc$	<b>0.082</b>	<b>0.075</b>	<b>0.196</b>	<b>0.077</b>	<b>0.052</b>	<b>0.125</b>	<b>0.103</b>	NA
	$\mathcal{D}_f$	1.000	ULM	0.235	0.150	0.239	0.291	0.340	0.324	0.267	NA
			RCM	0.312	0.567	0.352	0.364	0.478	0.538	0.429	NA
			$\Delta Acc$	<b>0.077</b>	<b>0.417</b>	<b>0.113</b>	<b>0.073</b>	<b>0.138</b>	<b>0.215</b>	<b>0.162</b>	NA

Table 2: Performance comparison of MRA (closed-source case) on various SOTA MU methods, where ULM indicates the  $Acc$  after MU while RCM indicates the  $Acc$  after MRA, and  $\Delta Acc$  shows the improvement (the Top-2  $\Delta Acc$  are marked in red while the Lowest-2 are marked in blue). TRM illustrates the  $Acc$  of original model.

sive evaluation. CIFAR-10 and CIFAR-100 [Krizhevsky *et al.*, 2009] consist of 60,000 low-resolution images classified into 10 and 100 classes, respectively. Oxford-IIIT Pet (Pet-37) [Parkhi *et al.*, 2012] contains 7,349 high-resolution images of cats and dogs in 37 classes. In particular, Oxford 102 Flower (Flower-102) [Nilsback and Zisserman, 2008] has 8,189 high-resolution images in 102 classes, where each class consists of between 40 and 258 images.

As shown in Table 1, we use the official splits of the training dataset  $\mathcal{D}_{tr}$  and the testing dataset  $\mathcal{D}_{ts}$  provided in the dataset package. For each dataset, five classes are selected to construct the forgetting dataset  $\mathcal{D}_f$  by randomly sampling 50% of data from  $\mathcal{D}_{tr}$ , and the prediction dataset used for evaluation is constructed by mixing  $\mathcal{D}_{ts}$  and  $\mathcal{D}_f$ , i.e.,  $\mathcal{D}_p = \mathcal{D}_{ts} \cup \mathcal{D}_f$  to jointly assess the prediction and recovery capability of RCM.

### Model Configuration

The proposed MRA framework is model-agnostic, so we use EfficientNet (EFN) [Tan and Le, 2019] on CIFAR datasets, ResNet [He *et al.*, 2016] on Pet-37, and Swin-Transformer (Swin-T) [Liu *et al.*, 2022] on Flower-102 for a comprehensive study. Moreover, in our framework, the STM does not necessarily have the same architecture as the ULM (Teacher Model). Therefore, we also evaluate the case of heterogeneous architectures on Flower-102, where the ULM is based

on Swin-T while the STM is based on ResNet, as shown in Tables 2 and 3.

We use SGD optimizer for MU methods with a momentum of 0.9, a weight decay of 0.005, and AdamW for our MRA scheme. Other more detailed settings can be found in the online extended version. For each comparison model, we carefully tuned their hyperparameters to achieve optimal performance.

### MU Methods for MRA Evaluation

In the experiments, a set of SOTA MU methods, including GA [Graves *et al.*, 2021], RL [Golatkhar *et al.*, 2020], FF [Becker and Liebig, 2022], IU [Koh and Liang, 2017], BU [Chen *et al.*, 2023], L1-SP [Jia *et al.*, 2023], UNSC [Chen *et al.*, 2024] and SalUn [Fan *et al.*, 2024], are involved to comprehensively evaluate the recall capability of proposed MRA. In particular, UNSC does not support the Swin-T architecture, so the evaluation results on Flower-102 are not available.

### 5.2 MRA in The Closed-source Case

First, we evaluated the performance of MRA in the closed-source case, i.e., ULMs are used as a black-box service where the model parameters are not accessible.

#### Overall MRA Efficacy Analysis

Table 2 show the accuracy ( $Acc$ ) of prediction dataset  $\mathcal{D}_p$  in terms of its subsets  $\mathcal{D}_{ts}$  and  $\mathcal{D}_{ts}$  respectively. We compared



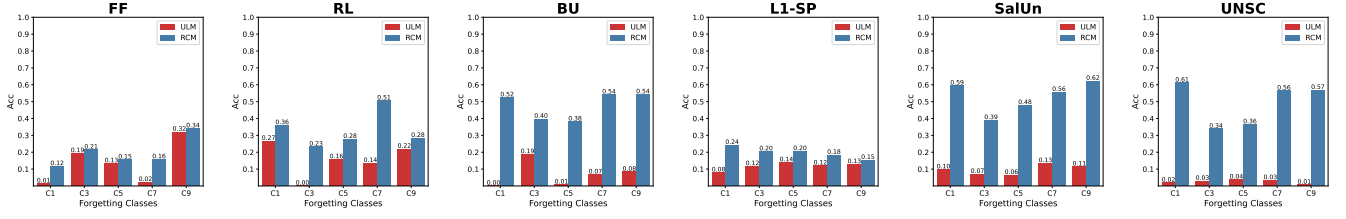


Figure 3: Comparison of the  $Acc$  between ULM and RCM (**closed-source case**) after MRA w.r.t. each forgetting class on CIFAR-10 dataset. Due to space limit, **LARGER** figures can be found in the online extended version.

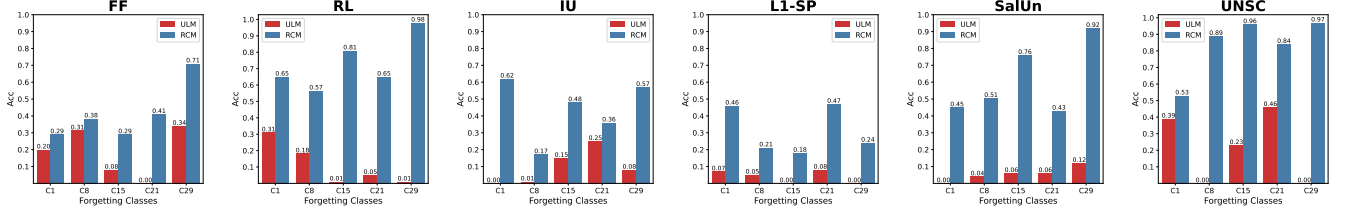


Figure 4: Comparison of the  $Acc$  between ULM and RCM (**closed-source case**) w.r.t. each forgetting class on Pet-37 dataset.

the  $Acc$  of ULM and RCM over four datasets and diverse configurations of the MU model. From the point of view of this table, all of the  $Acc$  on  $\mathcal{D}_f$  of TRM are 1.000 after training, while they drop to low  $Acc$  after MU, illustrating that all selected MU methods can successfully mitigate the influence of  $\mathcal{D}_f$  from well-trained models. Furthermore, according to all the results of improvement ( $\Delta Acc$ ) on both  $\mathcal{D}_f$  and  $\mathcal{D}_{ts}$ , we find that RCMs can unexceptionally improve the prediction accuracy on all datasets for all MU models, which overall proves that the proposed MRA is an effective and versatile model-agnostic framework to recover the class memberships of forgotten instances from ULMs.

More specifically, the improvement ( $\Delta Acc$ ) of **L1-SP** is overall smaller than that of other MU methods on both  $\mathcal{D}_f$  and  $\mathcal{D}_{ts}$ . This can be attributed to the weight pruning on TRM, which makes the parameters of ULM significantly different from those of TRM. As a result, the knowledge distillation from ULM is prone to having higher label noise. In contrast, the improvement ( $\Delta Acc$ ) of **SalUn** and **UNSC** is overall larger than that of other MU methods. This is because most MU methods degrade the model performance after unlearning, known as “over-unlearning”. In comparison, **SalUn** and **UNSC**, can precisely unlearn target forgetting samples without over-unlearning. As a result, **SalUn** and **UNSC** can provide less noisy pseudo labels for knowledge distillation, leading to better recall from ULMs.

According to the observation and analysis above, it reveals a phenomenon that “*The MU methods in precise unlearning may lead to high success rate to recall the forgotten class memberships via MRA*”. As a result, MRA can serve as a valuable tool to assess the potential risk of privacy leakage for MU methods, thus facilitating the development of more robust MU models.

### Demonstration of Class-specific Recovery Efficacy

To intuitively demonstrate the capacity of MRA to recall the forgotten class memberships on the prediction images  $\mathcal{X}_f$ , we further conducted detailed evaluations with respect to each

forgetting class. Figures 3 and 4 demonstrate the comparison of the  $Acc$  between ULM and RCM after MRA on CIFAR-10 and Pet-37 for each forgetting class. By checking the improvement of  $Acc$  for each class, we can observe a similar phenomenon as shown in Table 2. For example, the improvement of  $Acc$  for each class on **L1-SP** is relatively small due to over-unlearning. In comparison, the ULM via **SalUn** can precisely unlearn the target forgetting images, that is,  $\mathcal{X}_f$ , which leads to very low  $Acc$  for each forgetting class, whereas the improvement for each class after MRA is the most significant. That is, “*precise forgetting, easy recalling*”.

### 5.3 MRA in The Open-source Case

In comparison to the closed-source case, the ULMs are released with their parameters in the open-source case. As a result, the ULMs can be updated during the MRA process.

#### Overall MRA Efficacy Analysis

For the open-source case, we can find that the improvement ( $\Delta Acc$ ) of **SalUn** and **UNSC** is significant again due to the same reason presented above. Comparing Table 3 with Table 2, it is easy to find that the improvement of  $Acc$  in the open-source case method is significantly greater than that in the closed-source case. Especially, we find **IU** achieves the Lowest-2 improvement three times in the closed-source case but it achieves the Top-2 improvement three times in the open-source case. This is because the Confident Membership Recall step (cf. Algorithm 1) can effectively recall the knowledge retained by the ULM (as a noisy labeler) using confident pseudo-label samples. Then, the improved teacher model can distill less noisy knowledge to the STM. This alternative optimization process results in significant improvement.

#### Demonstration of Class-specific Recovery Efficacy

The analogy to the closed-source case, Figure 5 and appendix D.4 demonstrate the comparison of the  $Acc$  between ULM and RCM after MRA on CIFAR-10 and Pet-37 for each forgetting class in the open-source case. Compared to Figures 3 and 4, we can find that the gaps between different MU

		TRM		FF	RL	GA	IU	BU	L1-SP	SalUn	UNSC
CIFAR-10 T: EFN S: EFN	$\mathcal{D}_{ts}$	0.833	ULM	0.164	0.388	0.505	0.105	0.462	0.451	0.486	0.193
			RCM	0.708	0.471	0.582	0.689	0.783	0.479	0.829	0.774
			$\Delta Acc$	<b>0.544</b>	<b>0.083</b>	<b>0.077</b>	<b>0.584</b>	<b>0.321</b>	<b>0.027</b>	<b>0.343</b>	<b>0.581</b>
	$\mathcal{D}_f$	1.000	ULM	0.142	0.153	0.267	0.090	0.082	0.154	0.117	0.031
			RCM	0.726	0.351	0.509	0.759	0.877	0.338	0.997	0.921
			$\Delta Acc$	<b>0.584</b>	<b>0.197</b>	<b>0.242</b>	<b>0.668</b>	<b>0.795</b>	<b>0.184</b>	<b>0.880</b>	<b>0.890</b>
CIFAR-100 T: EFN S: EFN	$\mathcal{D}_{ts}$	0.643	ULM	0.159	0.593	0.531	0.140	0.529	0.537	0.410	0.275
			RCM	0.521	0.599	0.589	0.596	0.587	0.598	0.564	0.598
			$\Delta Acc$	<b>0.363</b>	<b>0.006</b>	<b>0.058</b>	<b>0.456</b>	<b>0.058</b>	<b>0.061</b>	<b>0.154</b>	<b>0.324</b>
	$\mathcal{D}_f$	1.000	ULM	0.094	0.086	0.201	0.201	0.254	0.199	0.227	0.209
			RCM	0.686	0.898	0.682	0.967	0.970	0.764	0.945	0.985
			$\Delta Acc$	<b>0.592</b>	<b>0.813</b>	<b>0.481</b>	<b>0.766</b>	<b>0.715</b>	<b>0.565</b>	<b>0.718</b>	<b>0.776</b>
Pet-37 T: ResNet S: ResNet	$\mathcal{D}_{ts}$	0.895	ULM	0.257	0.754	0.740	0.622	0.646	0.740	0.706	0.769
			RCM	0.782	0.850	0.843	0.838	0.846	0.841	0.849	0.856
			$\Delta Acc$	<b>0.525</b>	<b>0.096</b>	<b>0.103</b>	<b>0.216</b>	<b>0.200</b>	<b>0.102</b>	<b>0.143</b>	<b>0.087</b>
	$\mathcal{D}_f$	1.000	ULM	0.224	0.332	0.224	0.200	0.084	0.196	0.128	0.304
			RCM	0.904	0.960	0.952	0.936	0.960	0.920	0.948	0.976
			$\Delta Acc$	<b>0.680</b>	<b>0.628</b>	<b>0.728</b>	<b>0.736</b>	<b>0.876</b>	<b>0.724</b>	<b>0.820</b>	<b>0.672</b>
Flower102 T: Swin-T S: ResNet	$\mathcal{D}_{ts}$	0.939	ULM	0.294	0.756	0.314	0.512	0.706	0.530	0.496	NA
			RCM	0.607	0.915	0.772	0.829	0.889	0.856	0.857	NA
			$\Delta Acc$	<b>0.313</b>	<b>0.159</b>	<b>0.458</b>	<b>0.318</b>	<b>0.183</b>	<b>0.325</b>	<b>0.361</b>	NA
	$\mathcal{D}_f$	1.000	ULM	0.235	0.150	0.239	0.291	0.340	0.324	0.267	NA
			RCM	0.482	0.988	0.709	0.725	0.972	0.935	0.915	NA
			$\Delta Acc$	<b>0.247</b>	<b>0.838</b>	<b>0.470</b>	<b>0.433</b>	<b>0.632</b>	<b>0.611</b>	<b>0.648</b>	NA

Table 3: Performance comparison of MRA (**open-source case**) on various SOTA MU methods, where ULM indicates the  $Acc$  after MU while RCM indicates the  $Acc$  after MRA, and  $\Delta Acc$  shows the improvement (the Top-2  $\Delta Acc$  are marked in **red** while the Lowest-2 are marked in **blue**). TRM illustrates the  $Acc$  of original model.

methods are much smaller in Figure 5 and appendix D.4. Especially, we can find the improvement of  $Acc$  on different MU methods after the MRA is close for each forgetting class in Appendix D.4. Even the over-unlearning models, the ULMs via **IU** and **L1-SP** are effectively recalled their forgotten instances through the MRA process in terms of the balanced class membership recall strategy (cf. Section 4.3).

## 5.4 Ablation Study

In this section, we discuss the effectiveness of each key component in the implementation of the MRA framework. Since the MRA framework consists of two alternative learning steps, we will evaluate the following components.

**DST**: This component is the *Denosing Knowledge Distillation* step presented in Section 4.3, which aims to distill knowledge from noisy labeling teacher  $\mathcal{M}(\Theta_U)$  to STM  $\hat{\mathcal{M}}(\Theta_S)$ .

**STU**: The component serves as the class membership recall process for STM, as presented in the *Confident Membership Recall* step. That is, **DST+STU** is equivalent to the closed-source case of MRA.

**TCH**: The component serves as the class membership recall process for teacher models, as presented in the *Confident Membership Recall* step. That is, **DST+STU+TCH** is equivalent to the open-source case of MRA.

Component			FF		BU		SalUn		UNSC	
DST	STU	TCH	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$
✓			0.223	0.168	0.659	0.128	0.698	0.108	0.767	0.320
✓	✓		0.486	0.388	0.785	0.856	0.778	0.712	0.799	0.660
✓	✓	✓	0.782	0.904	0.846	0.960	0.849	0.948	0.856	0.976

Table 4: Ablation results ( $Acc$ ) of MRA on Pet-37 dataset

## Comparison Results

Table 4 reports the results of MRA on Pet-37, where the  $Acc$  on testing dataset  $\mathcal{D}_{ts}$  and forgetting dataset  $\mathcal{D}_f$  is reported. Additional results for ablation on other datasets can be found in the online extended version.

From the results, we can easily find that **DST** underperforms **DST+STU** and **DST+STU+TCH**. This is because the **DST** component only distills knowledge inferred from the prediction dataset  $\mathcal{D}_p$  from the ULM to the student model, where the distilled knowledge inevitably contains label noise from the ULM. As a result, the  $Acc$  of **DST** is close to that of **ULM**. In comparison, **DST+STU** and **DST+STU+TCH** are based on the alternative learning process with both distillation and recall steps (cf. Algorithm 1). In the recall step, samples with high-confidence agreement are extracted, which can effectively recall forgotten class memberships. The model with all components, i.e., **DST+STU+TCH**, achieves the best performance because the additional **TCH** component can further

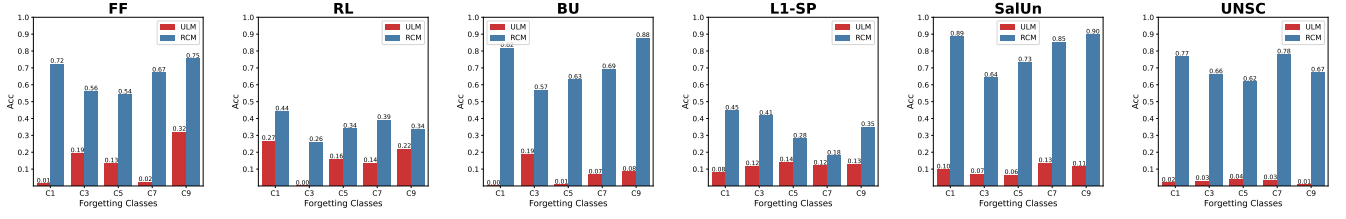


Figure 5: Comparison of the *Acc* between ULM and RCM (**open-source case**) after MRA w.r.t. each forgetting class on CIFAR-10 dataset. Due to space limit, **LARGER figures** can be found in the online extended version.

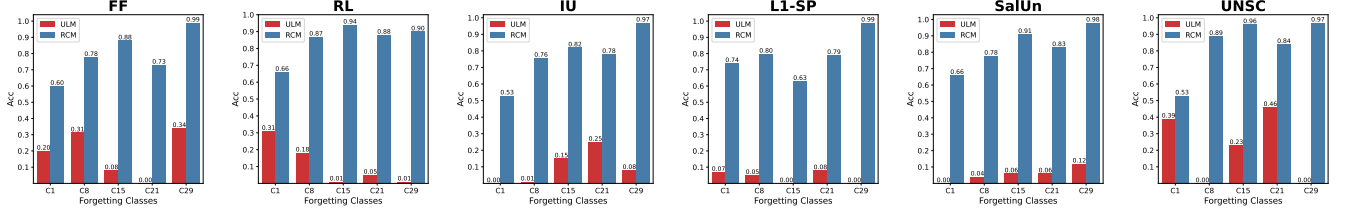


Figure 6: Comparison of the *Acc* between ULM and RCM (**open-source case**) w.r.t. each forgetting class on Pet-37 dataset.

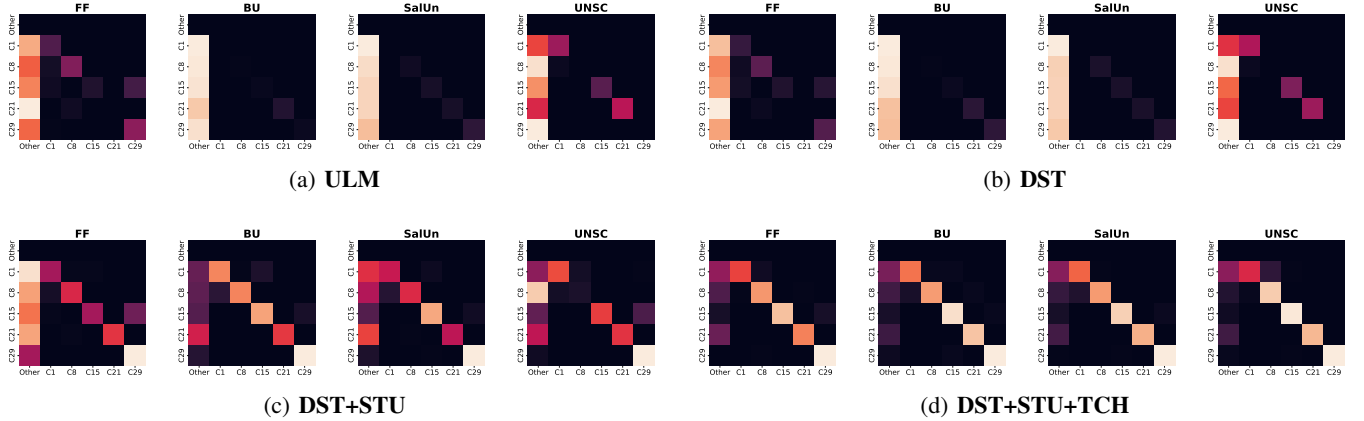


Figure 7: Confusion matrices with different configurations of MRA components. (a) demonstrates the confusion matrices after unlearning with different MU methods. Due to space limit, **LARGER figures** can be found in the online extended version.

recall the knowledge retained by the ULM, as illustrated in Section 5.3.

### Visualization of Confusion Matrices

Figure 7 (a-d) visualize the normalized confusion matrices on  $\mathcal{D}_f$  w.r.t. **ULM**, **DST**, **DST+STU** and **DST+STU+TCH**. From Figure 7 (a) **ULM**, we can observe the highlights in the first column of each subfigure, which illustrates that MU methods have successfully unlearned true class memberships of forgetting instances. From Figure 7 (b) to (c), the diagonals of confusion matrices w.r.t. **DST**, **DST+STU** and **DST+STU+TCH** become more and more noticeable, that is, more and more forgotten class memberships have been successfully recalled, which shows that each component plays an important role in the MRA framework.

## 6 Conclusion

This study is the first attempt to explore MRA against MU techniques to recall unlearned class memberships, highlight-

ing vulnerabilities of MU in data privacy protections. By using ULMs as noisy labelers, our implementation of MRA can recall forgotten class memberships from the ULMs without the need for the original model. Extensive experiments on four real-world datasets show that the proposed MRA framework exhibits high efficacy in recovering forgotten class memberships carried out by various MU methods. In particular, it reveals the phenomenon that “*The MU methods in precise unlearning may lead to high success rate to recall the forgotten class memberships via MRA*”. Consequently, the proposed MRA can serve as a valuable tool to assess the potential risk of privacy leakage for existing and new MU methods, thus gaining deeper insights into MU.

## Acknowledgments

This work is partially supported by the National Natural Science Foundation of China (NSFC Granted No. 62276190).



## References

- [Algan and Ulusoy, 2021] Gökrem Algan and Ilkay Ulusoy. Image classification with deep learning in the presence of noisy labels: A survey. *Knowledge-Based Systems*, 215:106771, 2021.
- [Bao et al., 2023] Guangyin Bao, Qi Zhang, Duoqian Miao, Zixuan Gong, and Liang Hu. Multimodal federated learning with missing modality via prototype mask and contrast. *CoRR*, abs/2312.13508, 2023.
- [Becker and Liebig, 2022] Alexander Becker and Thomas Liebig. Evaluating machine unlearning via epistemic uncertainty. *ArXiv preprint arXiv:2208.10836*, 2022.
- [Bourtoule et al., 2021] Lucas Bourtoule, Varun Chandrasekaran, Christopher A. Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *IEEE Symposium on Security and Privacy*, pages 141–159, 2021.
- [Cao and Yang, 2015] Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. *IEEE Symposium on Security and Privacy*, pages 463–480, 2015.
- [Carratino et al., 2022] Luigi Carratino, Moustapha Cissé, Rodolphe Jenatton, and Jean-Philippe Vert. On mixup regularization. *Journal of Machine Learning Research*, 23(325):1–31, 2022.
- [Chen et al., 2021] Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. When machine unlearning jeopardizes privacy. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, CCS ’21, page 896–911, New York, NY, USA, 2021. Association for Computing Machinery.
- [Chen et al., 2023] Min Chen, Weizhuo Gao, Gaoyang Liu, Kai Peng, and Chen Wang. Boundary unlearning: Rapid forgetting of deep networks via shifting the decision boundary. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7766–7775, 2023.
- [Chen et al., 2024] Huiqiang Chen, Tianqing Zhu, Xin Yu, and Wanlei Zhou. Machine unlearning via null space calibration. In Kate Larson, editor, *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24*, pages 358–366. International Joint Conferences on Artificial Intelligence Organization, 8 2024. Main Track.
- [Di et al., 2022] Jimmy Z. Di, Jack Douglas, Jayadev Acharya, Gautam Kamath, and Ayush Sekhari. Hidden poison: Machine unlearning enables camouflaged poisoning attacks. In *NeurIPS ML Safety Workshop*, 2022.
- [Fan et al., 2024] Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. In *The Twelfth International Conference on Learning Representations*, 2024.
- [Fredrikson et al., 2015] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In *ACM SIGSAC Conference on Computer and Communications Security*, pages 1322–1333, 2015.
- [Golatkar et al., 2020] Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9304–9312, 2020.
- [Graves et al., 2021] Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac machine learning. In *AAAI Conference on Artificial Intelligence*, volume 35, pages 11516–11524, 2021.
- [He et al., 2016] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [Hoofnagle et al., 2019] Chris Jay Hoofnagle, Bart van der Sloot, and Frederik J. Zuiderveen Borgesius. The european union general data protection regulation: what it is and what it means\*. *Information & Communications Technology Law*, 28:65 – 98, 2019.
- [Hu et al., 2023] Hongsheng Hu, Shuo Wang, Jiamin Chang, Haonan Zhong, Ruoxi Sun, Shuang Hao, Haojin Zhu, and Minhui Xue. A duty to forget, a right to be assured? exposing vulnerabilities in machine unlearning services. *ArXiv preprint arXiv:2309.08230*, 2023.
- [Hu et al., 2024] Hongsheng Hu, Shuo Wang, Tian Dong, and Minhui Xue. Learn what you want to unlearn: Unlearning inversion attacks against machine unlearning. In *IEEE Symposium on Security and Privacy*, 2024.
- [Itakura and Terada, 2018] Yoichiro Itakura and Mayu Terada. The significance and context of the establishment of california consumer privacy act of 2018. 2018.
- [Izzo et al., 2021] Zachary Izzo, Mary Anne Smart, Kamalika Chaudhuri, and James Zou. Approximate data deletion from machine learning models. In *International Conference on Artificial Intelligence and Statistics*, pages 2008–2016. PMLR, 2021.
- [Jia et al., 2023] Jinghan Jia, Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, Pranay Sharma, and Sijia Liu. Model sparsity can simplify machine unlearning. In *Neural Information Processing Systems*, 2023.
- [Koh and Liang, 2017] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International Conference on Machine Learning*, pages 1885–1894. PMLR, 2017.
- [Krizhevsky et al., 2009] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Master’s Thesis*, 2009.
- [Liu et al., 2022] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, Furu Wei, and Baining Guo. Swin transformer v2: Scaling up capacity and resolution. In *2022*

*IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11999–12009, 2022.

- [Lukasik *et al.*, 2020] Michal Lukasik, Srinadh Bhojanapalli, Aditya Menon, and Sanjiv Kumar. Does label smoothing mitigate label noise? In *International Conference on Machine Learning*, pages 6448–6458. PMLR, 2020.
- [Marchant *et al.*, 2022] Neil G. Marchant, Benjamin I. P. Rubinstein, and Scott Alfeld. Hard to forget: Poisoning attacks on certified machine unlearning. In *AAAI Conference on Artificial Intelligence*, volume 36, pages 7691–7700, 2022.
- [Miao *et al.*, 2024] Jiaxing Miao, Liang Hu, Qi Zhang, and Longbing Cao. Graph memory learning: Imitating lifelong remembering and forgetting of brain networks. *CoRR*, abs/2407.19183, 2024.
- [Nilsback and Zisserman, 2008] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *ICVGIP*, 2008.
- [Parkhi *et al.*, 2012] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *CVPR*, 2012.
- [Shokri *et al.*, 2017] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *IEEE Symposium on Security and Privacy*, pages 3–18, 2017.
- [Tan and Le, 2019] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- [Thudi *et al.*, 2022] Anvith Thudi, Gabriel Deza, Varun Chandrasekaran, and Nicolas Papernot. Unrolling sgd: Understanding factors influencing machine unlearning. In *IEEE European Symposium on Security and Privacy*, pages 303–319, 2022.
- [Wu *et al.*, 2024] Zhuojia Wu, Qi Zhang, Duoqian Miao, Kun Yi, Wei Fan, and Liang Hu. Hydiscgan: A hybrid distributed cgan for audio-visual privacy preservation in multimodal sentiment analysis. In *IJCAI*, pages 6550–6558. ijcai.org, 2024.
- [Zhang *et al.*, 2018] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *International Conference on Learning Representations*, 2018.

## Appendix

### A Source Code

To ensure a clear and concise code structure, we provide only the core implementation in the supplementary materials. Experimental datasets, trained model checkpoints, and other auxiliary files are excluded. The source code is available for review in the anonymous repository: <https://anonymous.4open.science/r/mra4mu-ijcai>.

### B Additional Experiment Details

This section provides additional information on **Data Preparation** and **Model Configuration** that are not specified in the main paper due to the limitation of space.

#### B.1 Dataset Preparation

For each dataset used in the experiments, five classes are selected to construct the forgetting dataset  $\mathcal{D}_f$ . Table 5 lists the class ID and name for each forgetting class.

#### B.2 Model Configuration

Tables 6 and 7 list the detailed model configurations of MRA across all the experiments for the closed-source and open-source cases respectively. The MRA model configuration can be separated into three groups: **Hyperparameter** and **Model Optimization**.

For **Network Architecture**, we use EfficientNet (EFN, small version) [Tan and Le, 2019] on CIFAR datasets, ResNet (ResNet18) [He *et al.*, 2016] on Pet-37, and Swin-Transformer V2 (Swin-T, tiny version) [Liu *et al.*, 2022] on Flower-102 for a comprehensive study. In particular, the ULM is based on Swin-T while the STM is based on ResNet for the Flower-102 dataset.

For **Hyperparameter**, it is further categorized into three subgroups, (1) **Mixup**:  $\alpha_x$  is the beta distribution hyperparameter for Mixup on images (cf. Eq. (1) in the main paper), and  $\alpha_y$  is the hyperparameter of beta distribution for the Mixup on labels (cf. Eq. (2) in the main paper); (2) **Smoothing**:  $\gamma_l$  is the hyperparameter of Laplace smoothing (cf. Eq. (6) in the main paper), and  $\gamma_s$  is the hyperparameter of label smoothing (cf. Eq. (8) in the main paper); (3) **Top-K**:  $\tau$  is the hyperparameter to control how many confidence-agreement samples are selected (cf. Eq. (8) in the main paper).

For **Model Optimization**, we use AdamW as the optimizer across all experiments,  $LR_S$  denotes the initial learning rate for STM, and  $LR_U$  denotes the initial learning rate for ULM, and the batch size for each dataset is also reported.

### C Additional Results of MRA Efficacy

To intuitively demonstrate the capacity of MRA to recall the forgotten class memberships on the prediction images  $\mathcal{X}_f$ , we conducted in-detail evaluations w.r.t. each forgetting class by comparing the  $Acc$  between ULM and RCM after MRA on CIFAR-10, CIFAR-100, Pet-37 and Flower-102.

From the illustrations in both the closed-source and the open-source case, we can clearly observe the phenomenon that “An MU method with more precise unlearning may lead to higher success rate of MRA to recall the forgotten class

memberships”, e.g., SalUn and UNSC. As a result, MRA can serve as a valuable tool to assess the potential risk of privacy leakage for MU methods, thus facilitating the development of more robust MU models.

#### C.1 Class-specific Recovery Efficacy (Close-source Case)

#### C.2 Class-specific Recovery Efficacy (Open-source Case)

### D Additional Results of Ablation Study

In this section, we demonstrate more results of the ablation study on CIFAR-10, CIFAR-100, Pet-37, and Flower-102. The ablation components have been presented in the main paper, that is:

- **DST**: This component is the *Denosing Knowledge Distillation* step presented in the main paper, which aims to distill knowledge from the teacher model to the student model.
- **STU**: The component serves as the class membership recall process for STM, as presented in the *Confident Membership Recall* step. That is, **DST+STU** is equivalent to the closed-source case of MRA.
- **TCH**: The component serves as the recall of class membership for teacher models as presented in the *Confident Membership Recall* step. That is, **DST+STU+TCH** is equivalent to the open-source case of MRA.

#### D.1 Ablation Study on CIFAR-10

#### D.2 Ablation Study on CIFAR-100

#### D.3 Ablation Study on Pet-37

#### D.4 Ablation Study on FLower-102

Dataset	Forgetting Classes ID and Name
CIFAR-10	C1: automobile; C3: cat; C5: dog; C7: horse; C9: truck
CIFAR-100	C10: bowl; C30: dolphin; C50: mouse; C70: rose; C90: train
Pet-37	C1: Bengal; C8: Ragdoll; C15: Basset Hound C21: Beagle; C29: Japanese Chin
Flower-102	C50: petunia; C72: water lily; C76: passion flower; C88: watercress; C93: foxglove;

Table 5: Details of five forgetting classes for each dataset

Dataset	Network Architecture		Hyperparameter					Model Optimization			
			Mixup		Smoothing		Top-K				
	ULM	STM	$\alpha_x$	$\alpha_y$	$\gamma_l$	$\gamma_s$	$\tau$	Optimizer	$LR_S$	$LR_U$	Batch Size
CIFAR-10	EFN	EFN	0.20	0.75	0.01	0.05	0.60	AdamW	1.00E-04	NA	256
CIFAR-100	EFN	EFN	0.20	0.75	0.01	0.05	0.60	AdamW	2.00E-04	NA	256
Pet-37	ResNet	ResNet	0.20	0.75	0.01	0.05	0.05	AdamW	2.00E-04	NA	64
Flower-102	Swin-T	ResNet	0.20	0.75	0.01	0.05	0.05	AdamW	2.00E-04	NA	32

Table 6: Detailed model configuration for MRA in the closed-source case (ULM is not updated so  $LR_U$  is NA)

Dataset	Network Architecture		Hyperparameter					Model Optimization			
			Mixup		Smoothing		Top-K				
	ULM	STM	$\alpha_x$	$\alpha_y$	$\gamma_l$	$\gamma_s$	$\tau$	Optimizer	$LR_S$	$LR_U$	Batch Size
CIFAR-10	EFN	EFN	0.20	0.75	0.01	0.05	0.05	AdamW	1.00E-04	5.00E-05	256
CIFAR-100	EFN	EFN	0.20	0.75	0.01	0.05	0.05	AdamW	2.00E-04	1.00E-04	256
Pet-37	ResNet	ResNet	0.20	0.75	0.01	0.05	0.05	AdamW	2.00E-04	1.00E-04	64
Flower-102	Swin-T	ResNet	0.20	0.75	0.01	0.05	0.05	AdamW	1.00E-04	1.00E-04	32

Table 7: Detailed model configuration for MRA in the open-source case (ULM is optimized with  $LR_U$ )

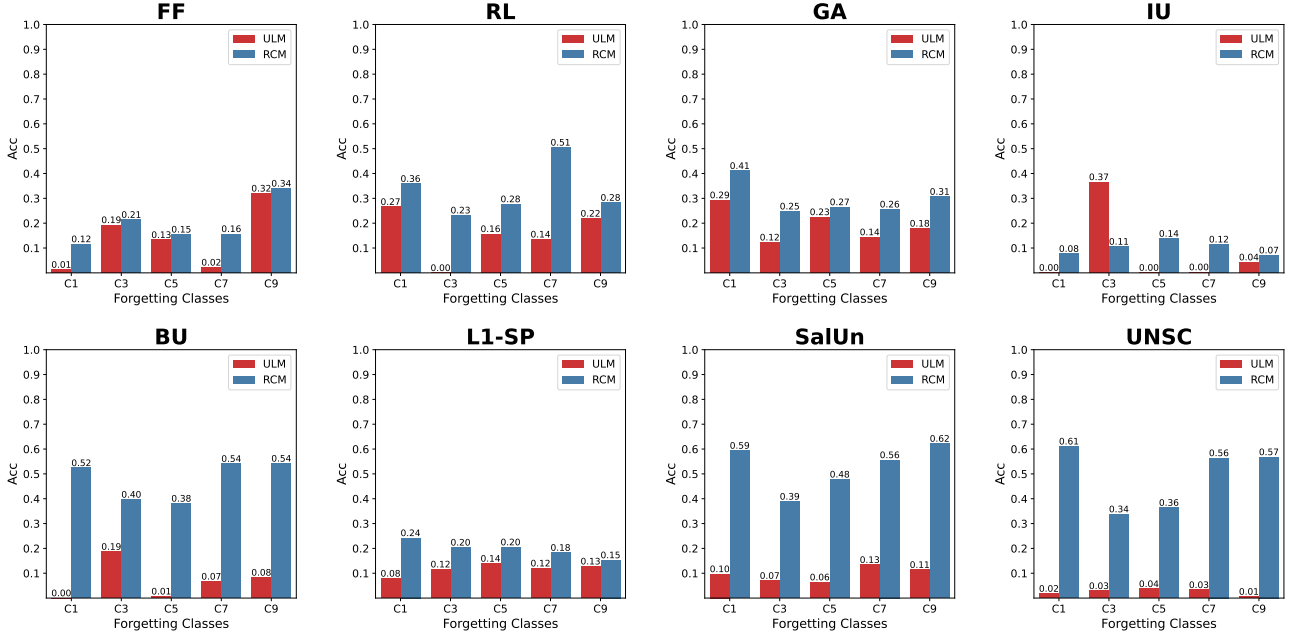


Figure 8: Comparison of the Acc between ULM and RCM (closed-source case) w.r.t. each forgetting class on CIFAR-10 dataset.

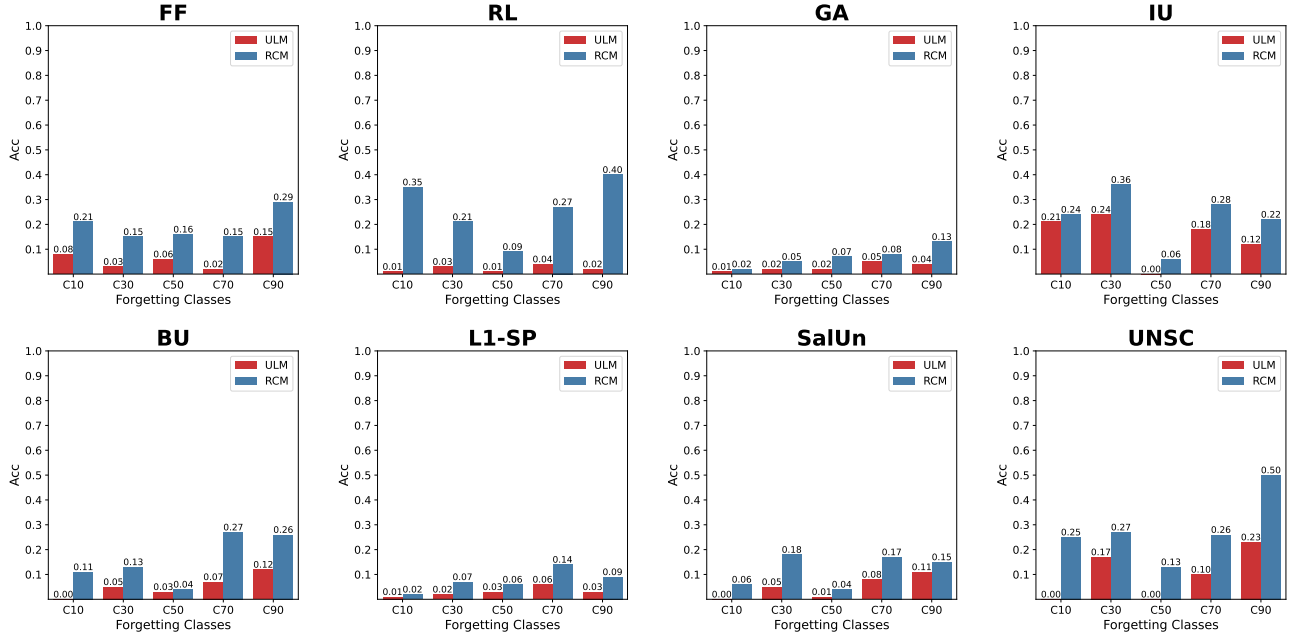


Figure 9: Comparison of the *Acc* between ULM and RCM (closed-source case) w.r.t. each forgetting class on CIFAR-100 dataset.

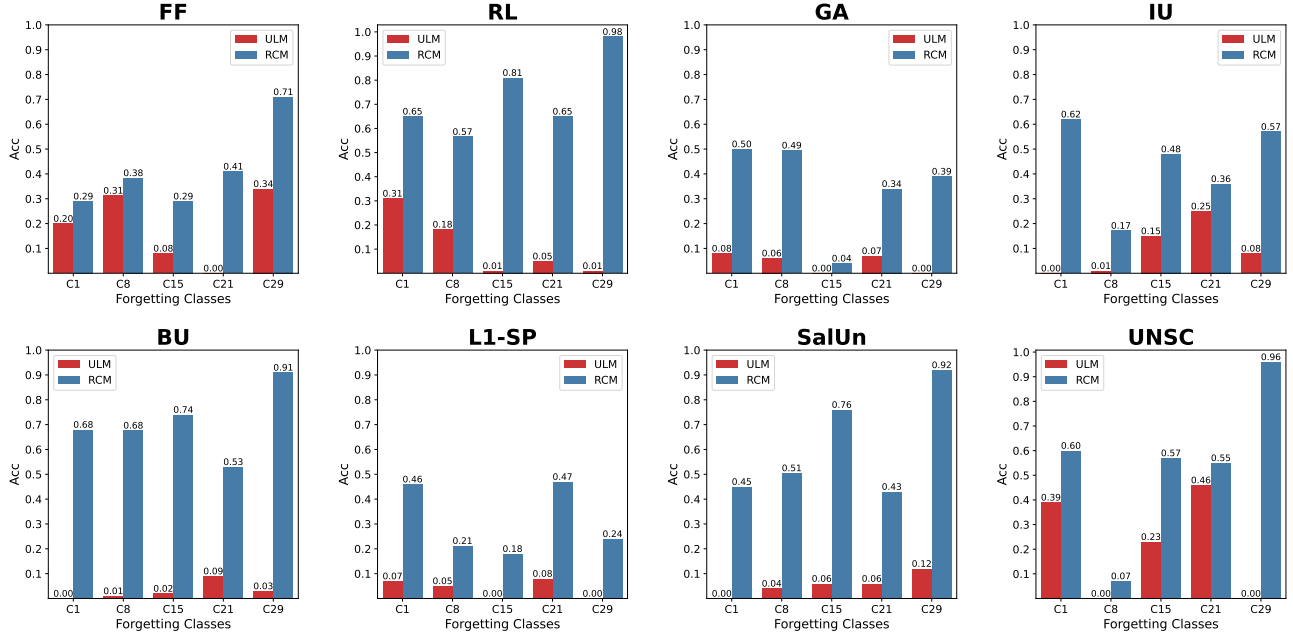


Figure 10: Comparison of the *Acc* between ULM and RCM (closed-source case) w.r.t. each forgetting class on Pet-37 dataset.

Component			FF		RL		GA		IU		BU		L1-SP		SalUn		UNSC	
DST	STU	TCH	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$
✓			0.163	0.126	0.393	0.144	0.501	0.252	0.114	0.101	0.458	0.059	0.448	0.131	0.451	0.027	0.199	0.030
✓	✓		0.252	0.195	0.479	0.328	0.528	0.358	0.115	0.096	0.651	0.507	0.456	0.231	0.697	0.615	0.568	0.513
✓	✓	✓	0.708	0.726	0.471	0.351	0.582	0.509	0.689	0.759	0.783	0.877	0.479	0.338	0.829	0.997	0.774	0.921

Table 8: Ablation results (*Acc*) of MRA on CIFAR-10 dataset w.r.t. different MU methods



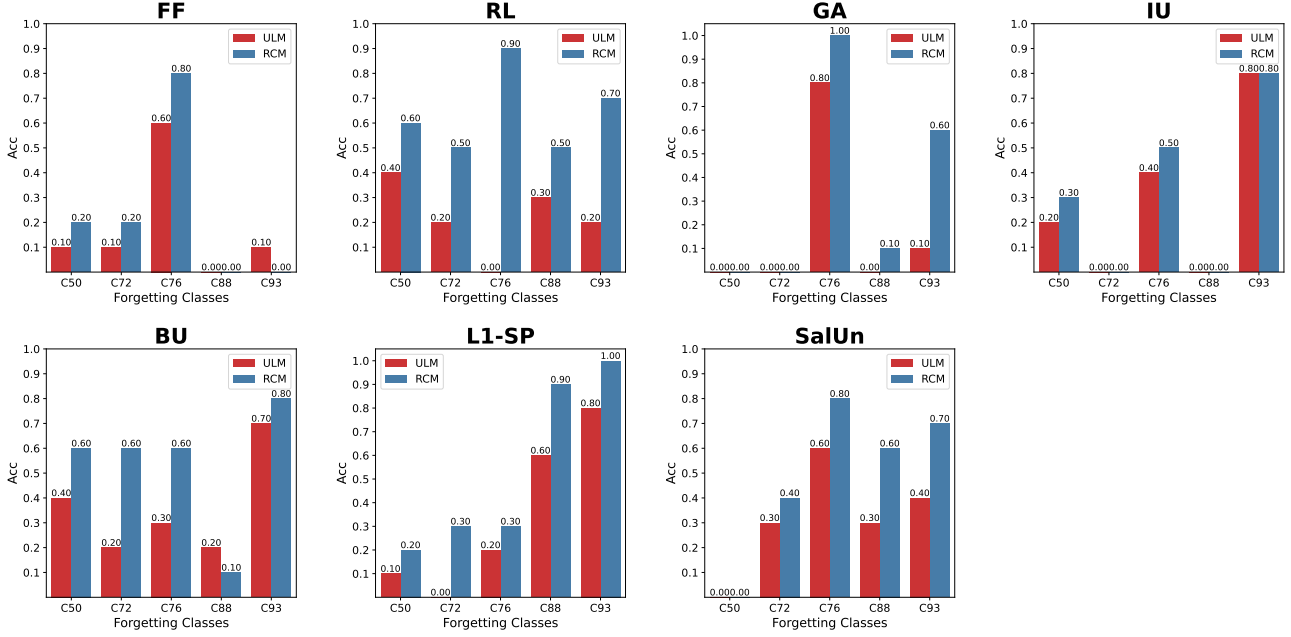


Figure 11: Comparison of the Acc between ULM and RCM (closed-source case) w.r.t. each forgetting class on Flower-102 dataset.

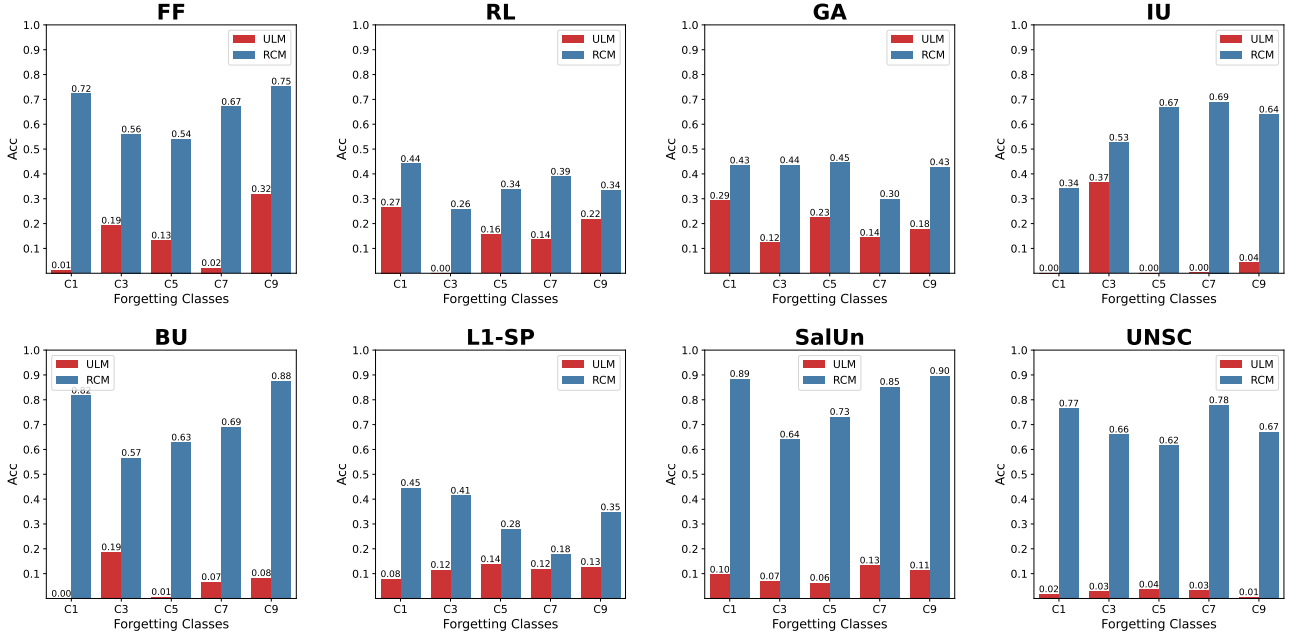


Figure 12: Comparison of the Acc between ULM and RCM (open-source case) w.r.t. each forgetting class on CIFAR-10 dataset.

Component			FF		RL		GA		IU		BU		L1-SP		SalUn		UNSC	
DST	STU	TCH	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$
✓			0.173	0.103	0.594	0.015	0.528	0.180	0.135	0.194	0.532	0.210	0.536	0.183	0.440	0.142	0.244	0.131
✓	✓		0.295	0.214	0.607	0.537	0.551	0.270	0.297	0.256	0.566	0.423	0.550	0.258	0.522	0.358	0.480	0.372
✓	✓	✓	0.521	0.686	0.599	0.898	0.589	0.682	0.596	0.967	0.587	0.970	0.598	0.764	0.564	0.945	0.598	0.985

Table 9: Ablation results (Acc) of MRA on CIFAR-100 dataset w.r.t. different MU methods

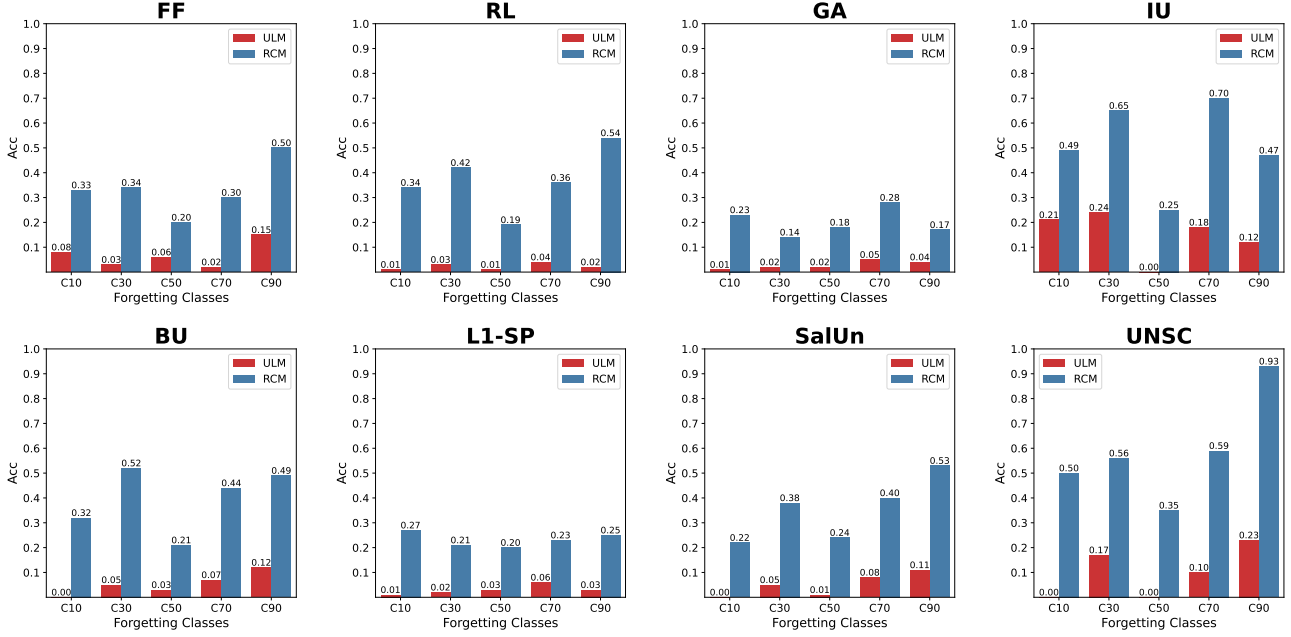


Figure 13: Comparison of the *Acc* between ULM and RCM (**open-source case**) w.r.t. each forgetting class on CIFAR-100 dataset.

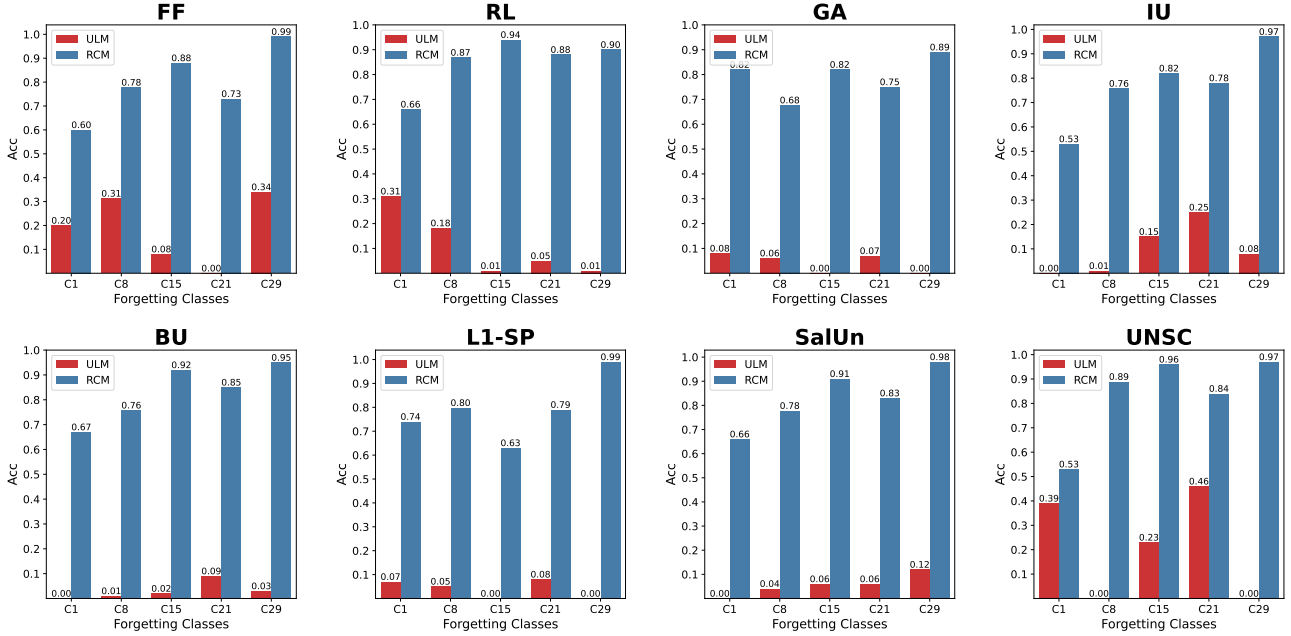


Figure 14: Comparison of the *Acc* between ULM and RCM (**open-source case**) w.r.t. each forgetting class on Pet-37 dataset.

Component			FF		RL		GA		IU		BU		L1-SP		SalUn		UNSC	
DST	STU	TCH	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$
✓			0.223	0.168	0.766	0.364	0.739	0.176	0.576	0.144	0.659	0.128	0.732	0.208	0.698	0.108	0.767	0.320
✓	✓		0.486	0.388	0.816	0.884	0.757	0.520	0.682	0.448	0.785	0.856	0.758	0.440	0.778	0.712	0.799	0.660
✓	✓	✓	0.782	0.904	0.850	0.960	0.843	0.952	0.838	0.936	0.846	0.960	0.841	0.920	0.849	0.948	0.856	0.976

Table 10: Ablation results (*Acc*) of MRA on Pet-37 dataset w.r.t. different MU methods

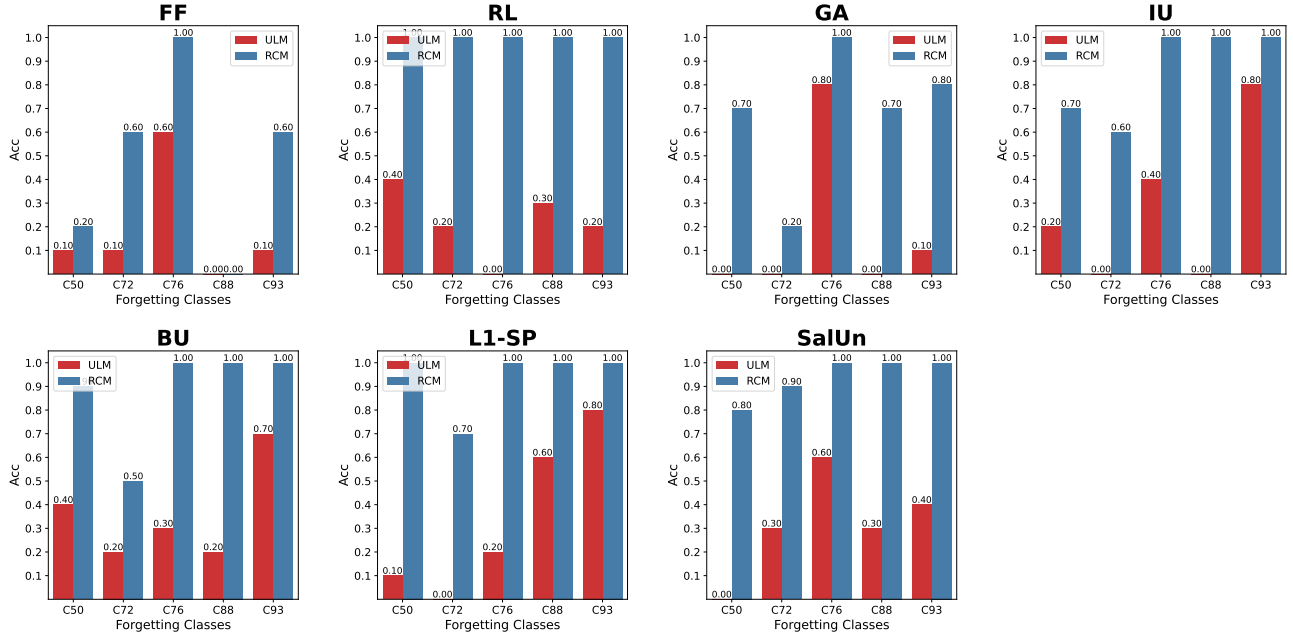


Figure 15: Comparison of the Acc between ULM and RCM (open-source case) w.r.t. each forgetting class on Flower-102 dataset.

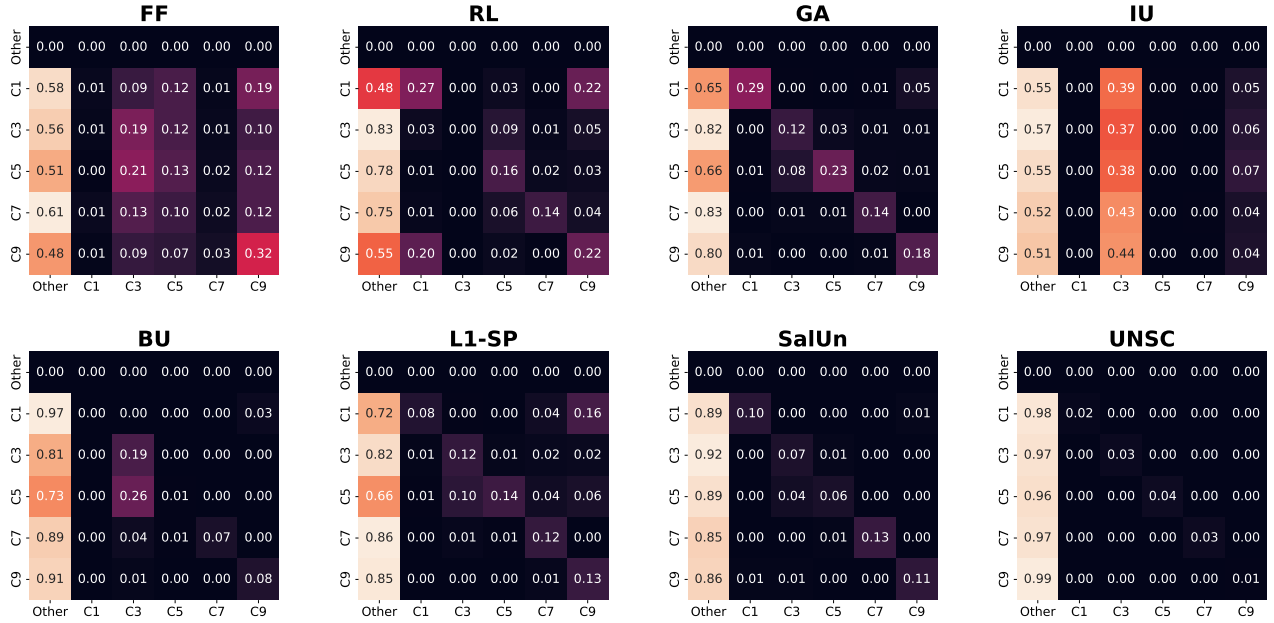


Figure 16: Confusion matrices of ULM w.r.t. different MU methods

Component			FF		RL		GA		IU		BU		L1-SP		SalUn		UNSC	
DST	STU	TCH	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$	$\mathcal{D}_{ts}$	$\mathcal{D}_f$
✓			0.278	0.198	0.741	0.170	0.319	0.231	0.519	0.320	0.714	0.377	0.511	0.308	0.505	0.247	NA	NA
✓	✓		0.376	0.312	0.831	0.567	0.510	0.352	0.589	0.364	0.758	0.478	0.655	0.538	0.599	0.429	NA	NA
✓	✓	✓	0.607	0.482	0.915	0.988	0.772	0.709	0.829	0.725	0.889	0.972	0.856	0.935	0.857	0.915	NA	NA

Table 11: Ablation results (Acc) of MRA on FLower-102 dataset w.r.t. different MU methods

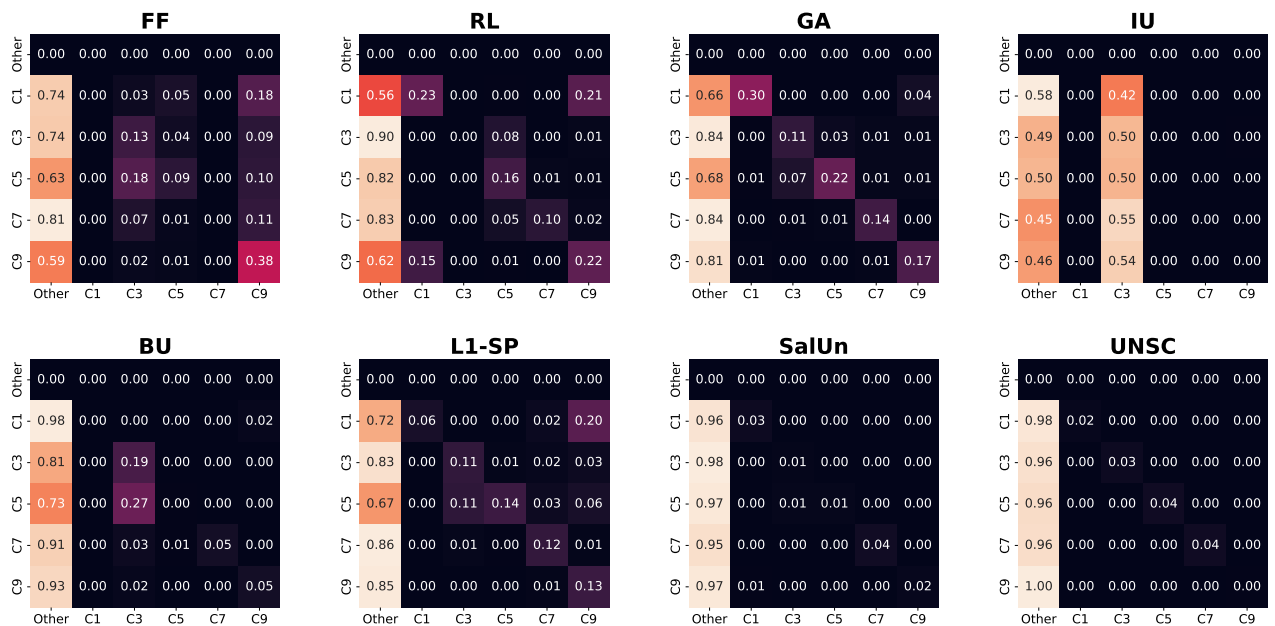


Figure 17: Confusion matrices of **DST** w.r.t. different MU methods

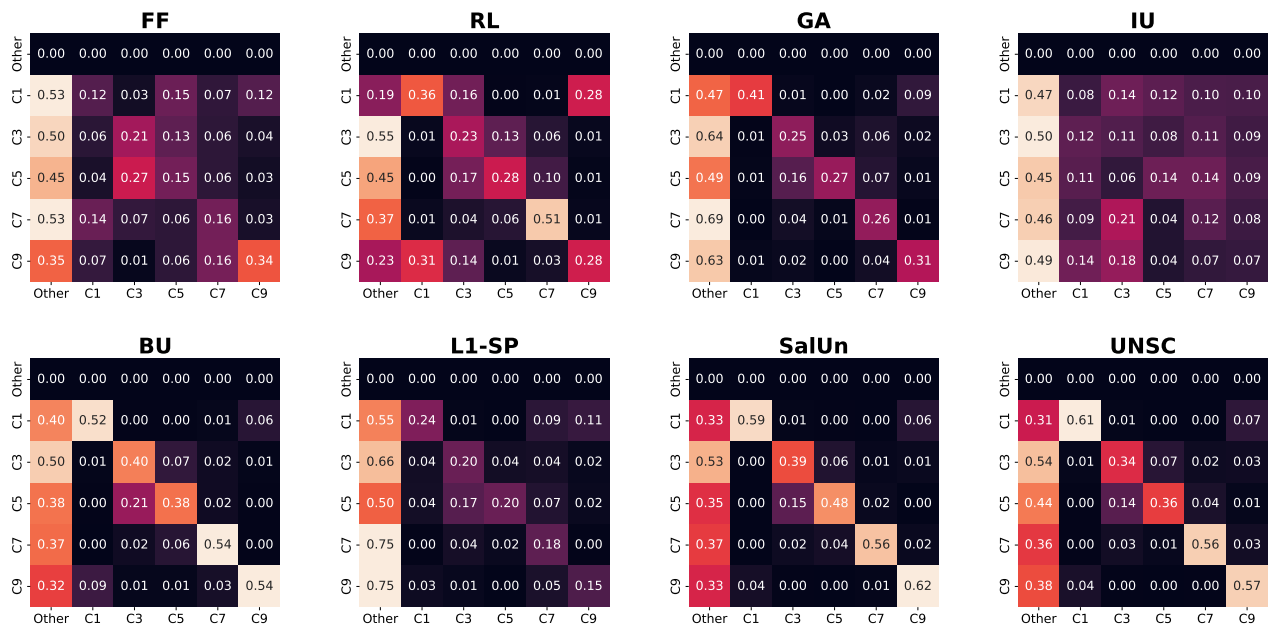


Figure 18: Confusion matrices of **DST+STU** w.r.t. different MU methods

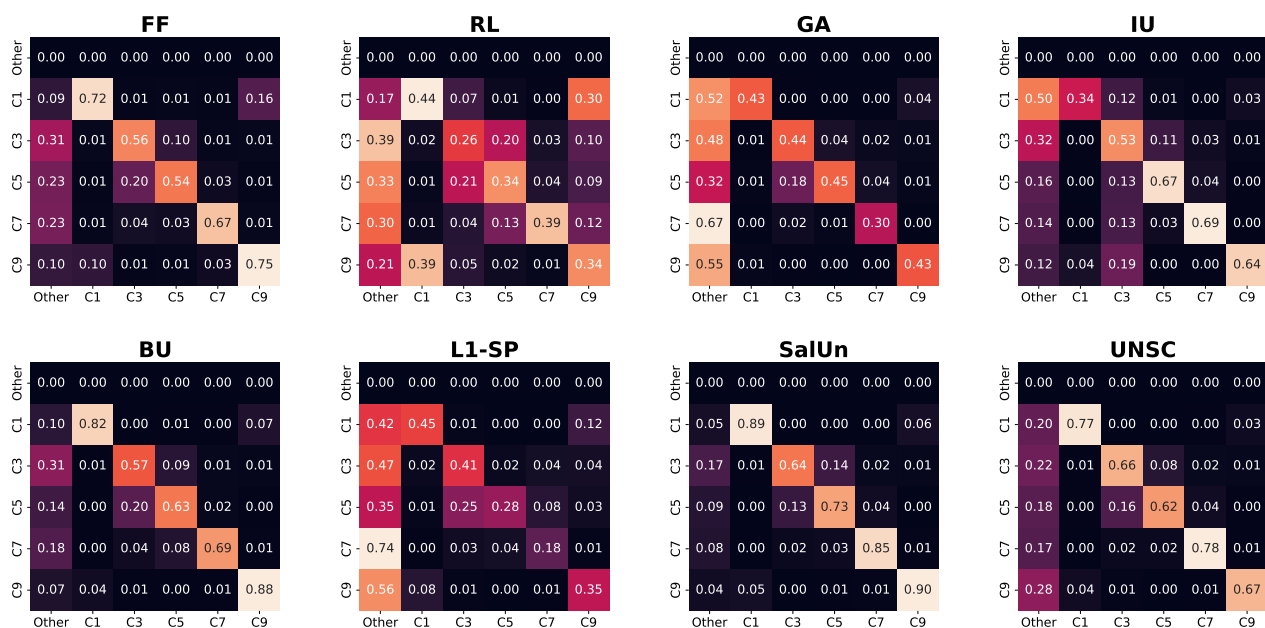


Figure 19: Confusion matrices of the DST+STU+TCH w.r.t. different MU methods

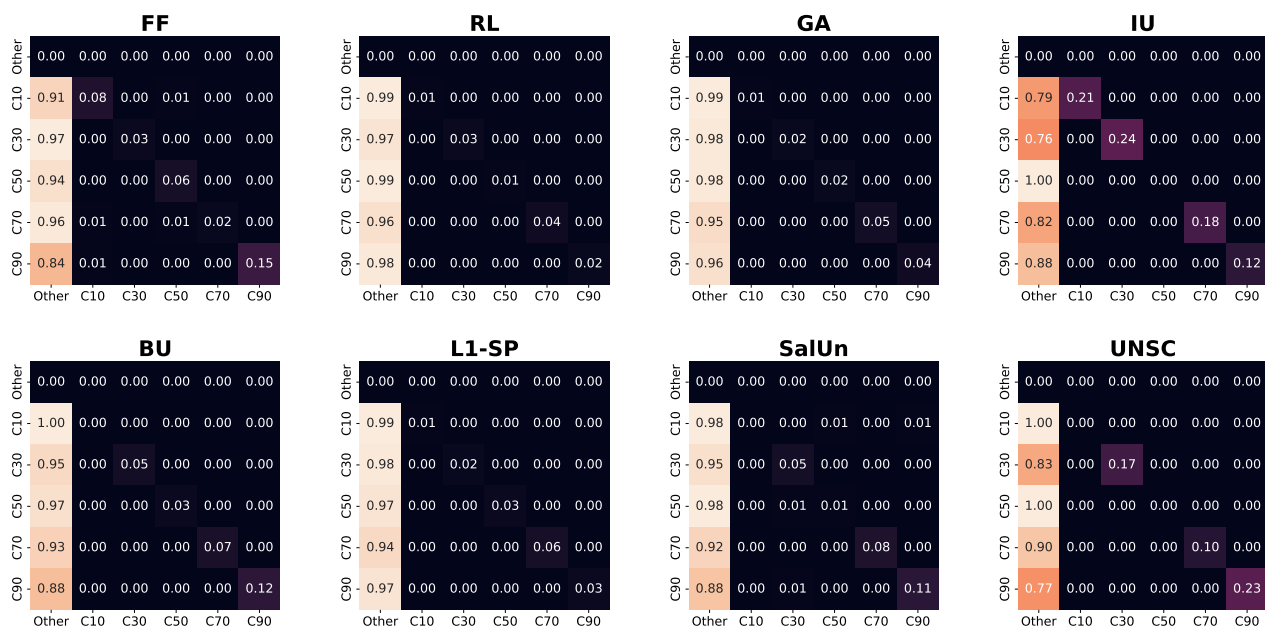


Figure 20: Confusion matrices of the UML w.r.t. different MU methods



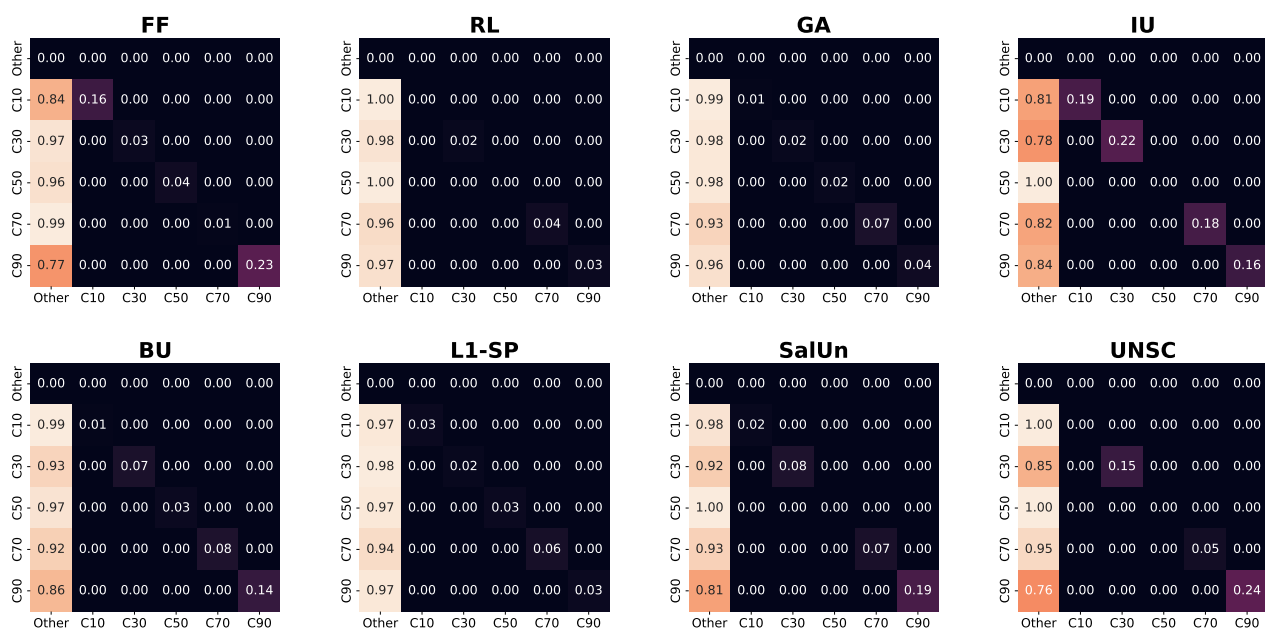


Figure 21: Confusion matrices of the **DST** w.r.t. different MU methods

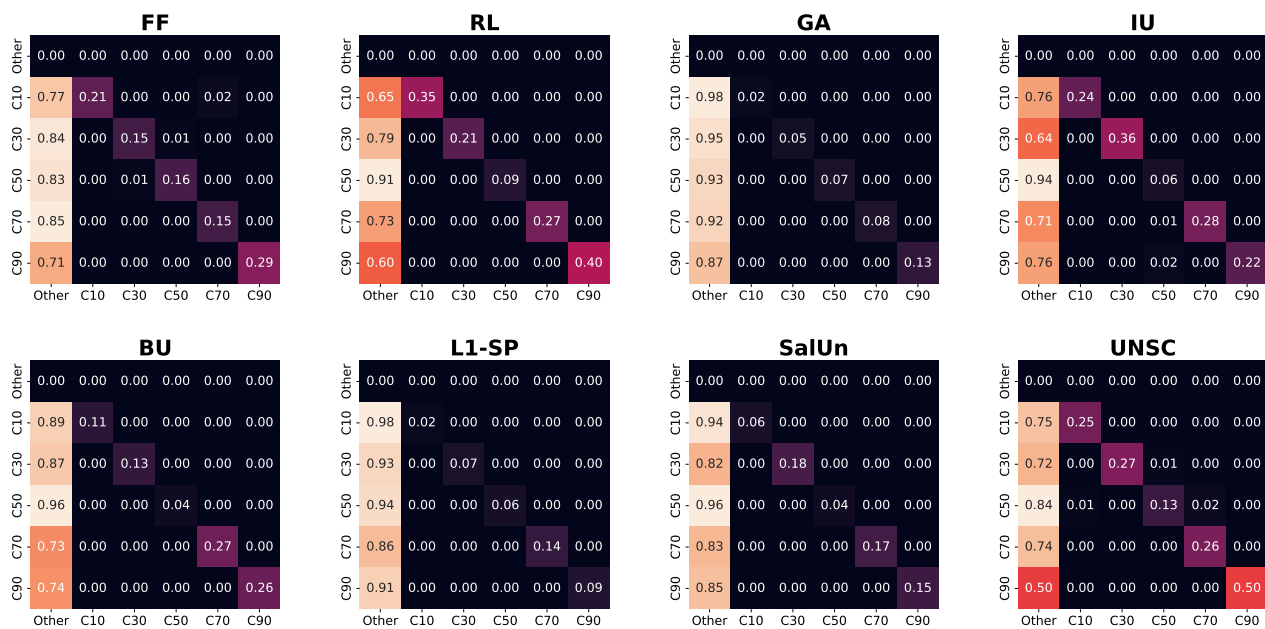


Figure 22: Confusion matrices of the **DST+STU** w.r.t. different MU methods

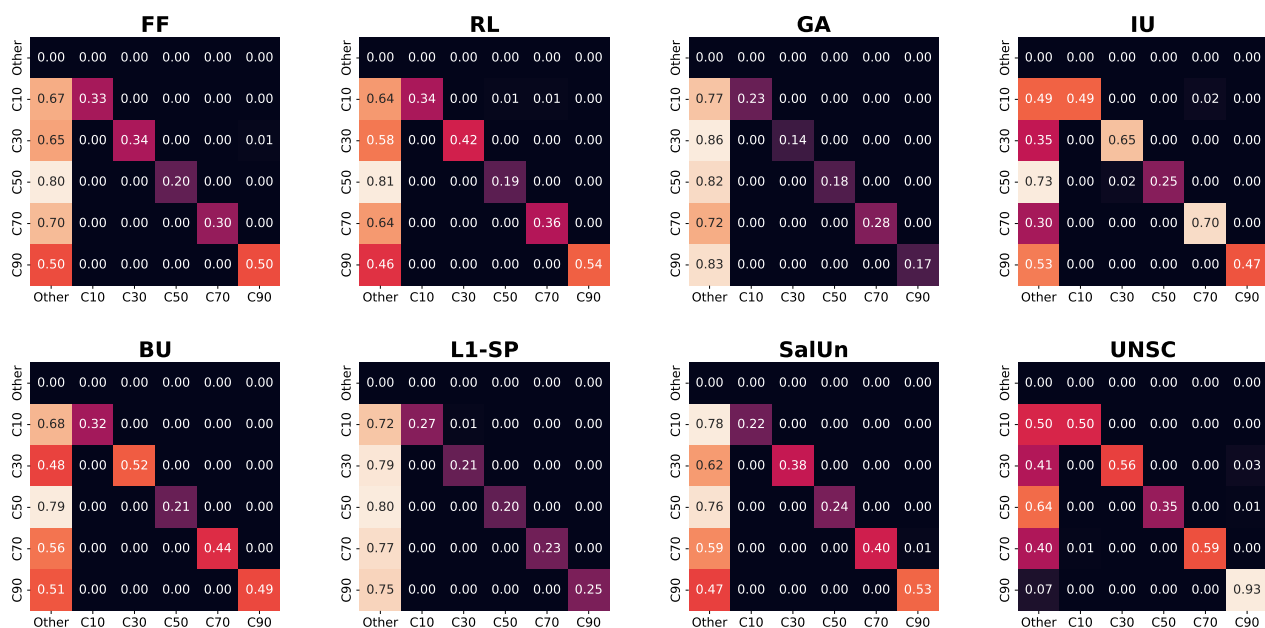


Figure 23: Confusion matrices of the DST+STU+TCH w.r.t. different MU methods

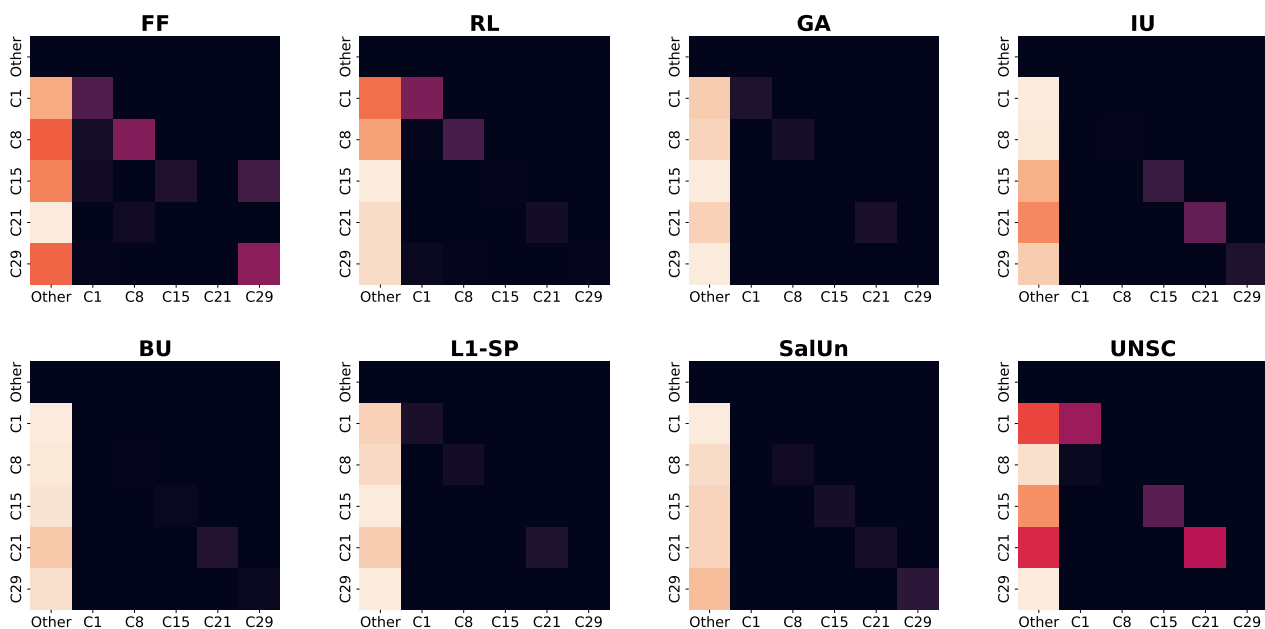


Figure 24: Confusion matrices of the UML w.r.t. different MU methods

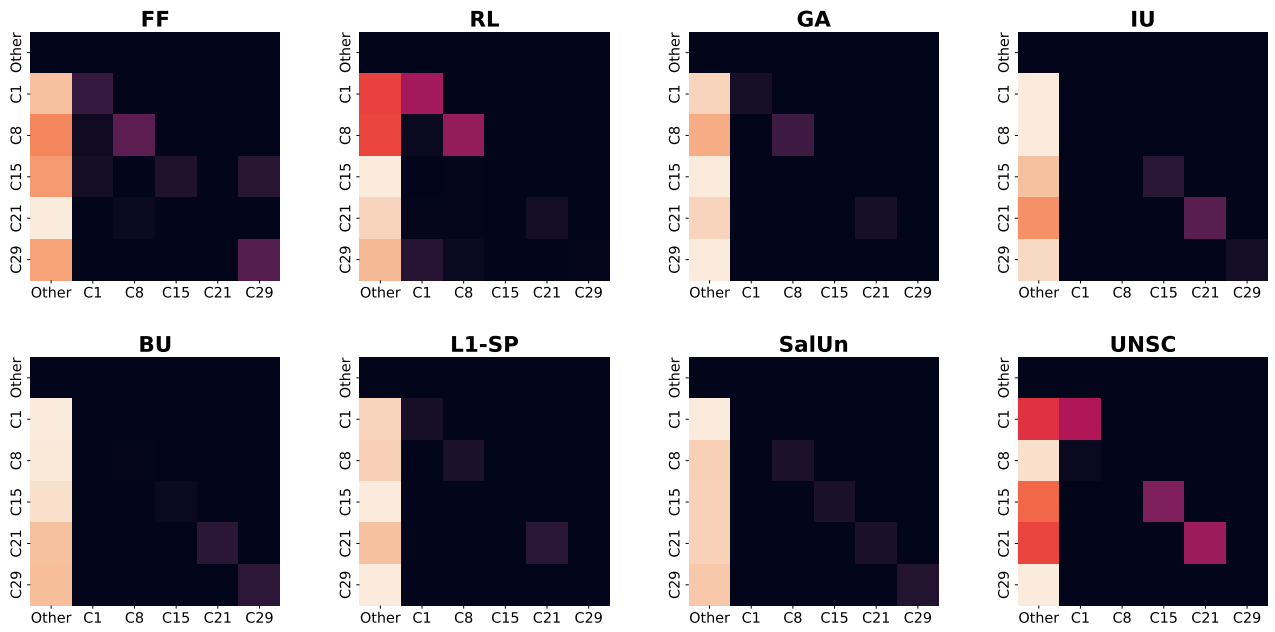


Figure 25: Confusion matrices of the **DST** w.r.t. different MU methods

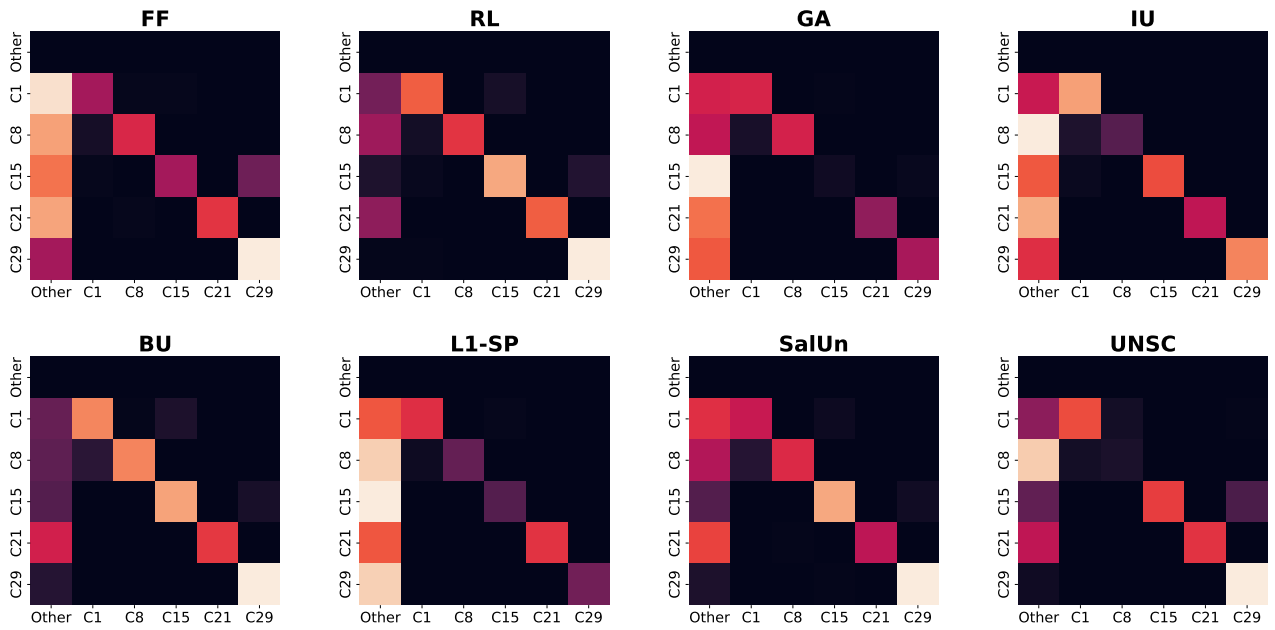


Figure 26: Confusion matrices of the **DST+STU** w.r.t. different MU methods

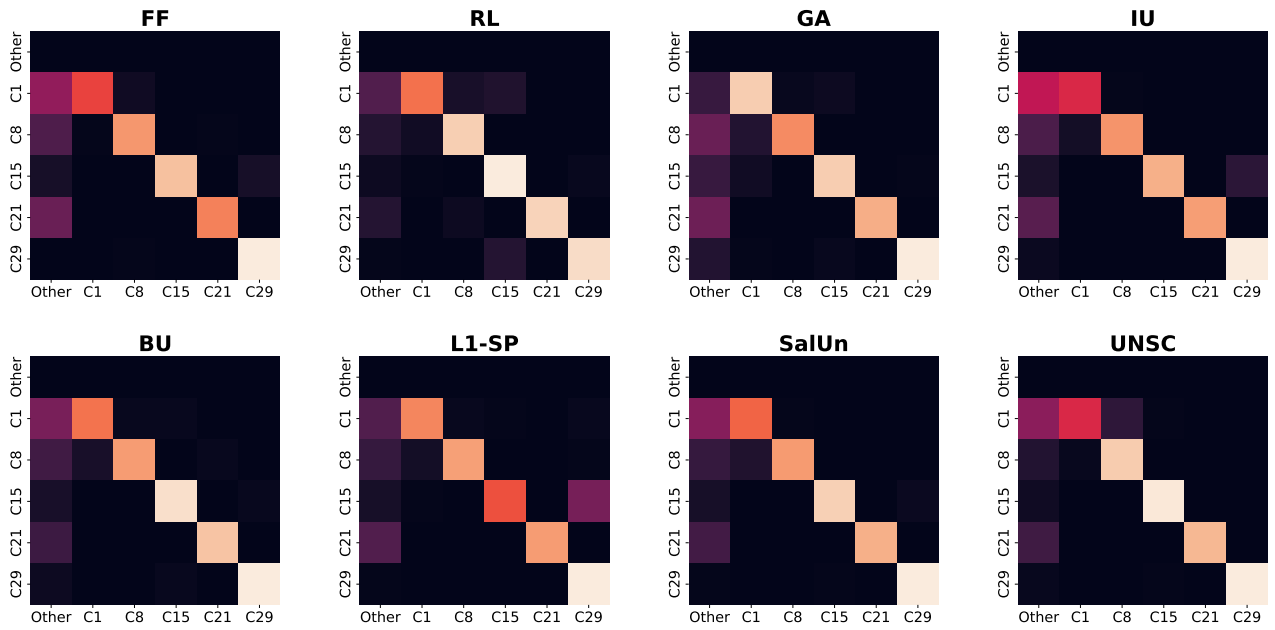


Figure 27: Confusion matrices of the **DST+STU+TCH** w.r.t. different MU methods

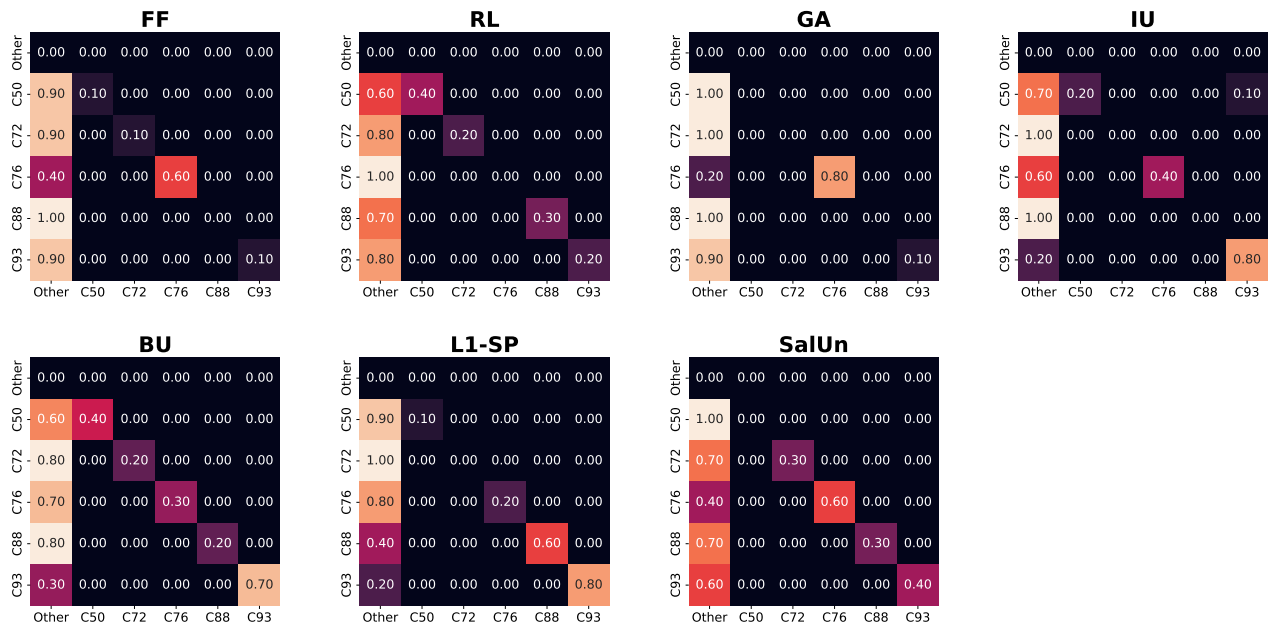


Figure 28: Confusion matrices of the **ULM** w.r.t. different MU methods

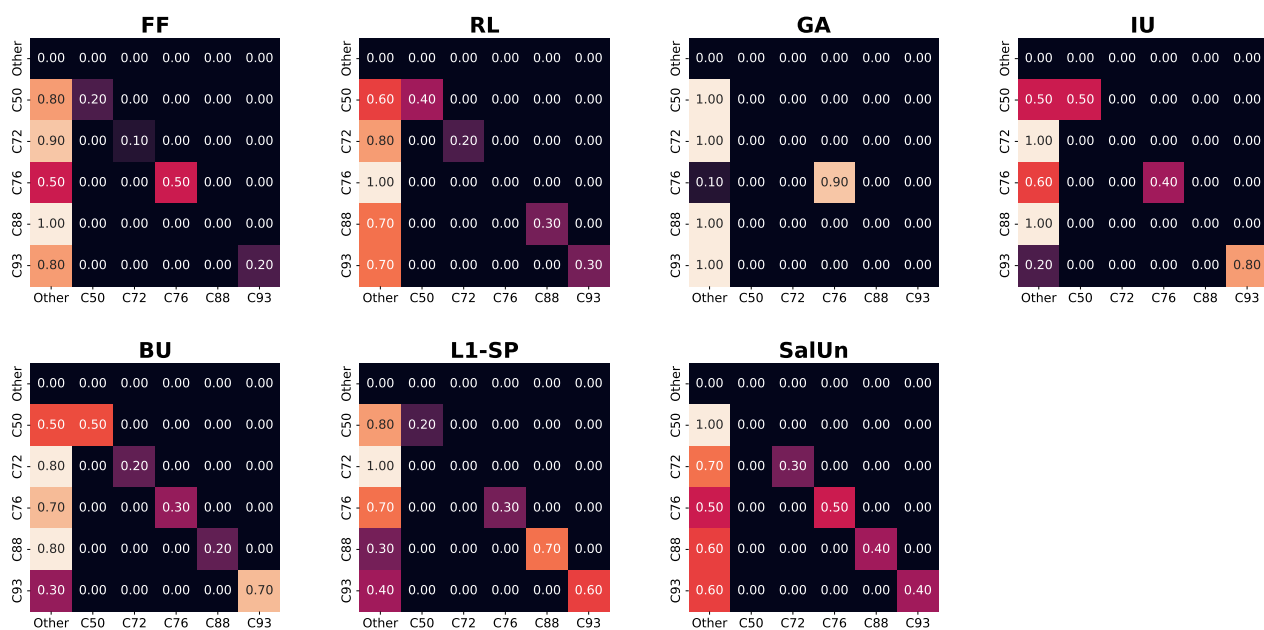


Figure 29: Confusion matrices of the **DST** w.r.t. different MU methods

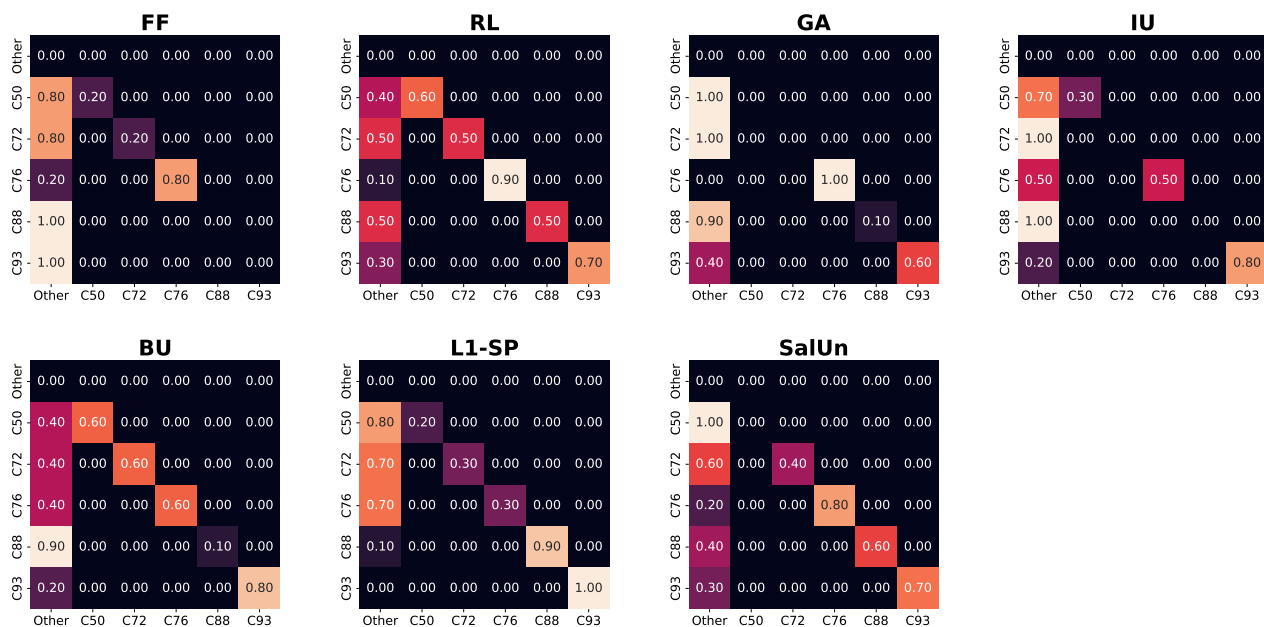


Figure 30: Confusion matrices of the **DST+STU** w.r.t. different MU methods



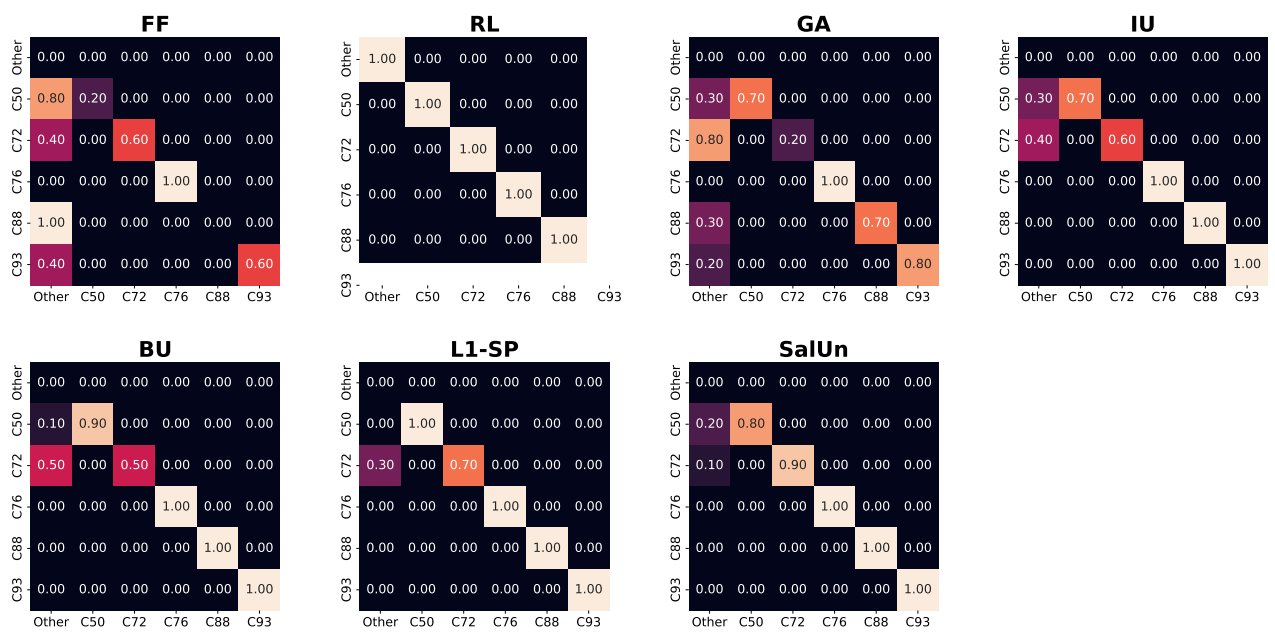


Figure 31: Confusion matrices of the **DST+STU+TCH** w.r.t. different MU methods