DECOUPLING THE CLASS LABEL AND THE TARGET CONCEPT IN MACHINE UNLEARNING

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026 027 028

029

Paper under double-blind review

ABSTRACT

Machine unlearning as an emerging research topic for data regulations, aims to adjust a trained model to approximate a retrained one that excludes a portion of training data. Previous studies showed that class-wise unlearning is effective in forgetting the knowledge of a training class, either through gradient ascent on the forgetting data or fine-tuning with the remaining data. However, while these methods are useful, they are insufficient as the class label and the target concept are often considered to coincide. In this work, we expand the scope by considering the label domain mismatch and investigate three problems beyond the conventional all matched forgetting, e.g., target mismatch, model mismatch, and data mismatch forgetting. We systematically analyze the new challenges in restrictively forgetting the target concept and also reveal crucial forgetting dynamics in the representation level to realize these tasks. Based on that, we propose a general framework, namely, TARget-aware Forgetting (TARF). It enables the additional tasks to actively forget the target concept while maintaining the rest part, by simultaneously conducting annealed gradient ascent on the forgetting data and selected gradient descent on the hard-to-affect remaining data. Empirically, various experiments under our newly introduced settings are conducted to demonstrate the effectiveness of our TARF.

1 INTRODUCTION

In response to data regulations [26, 49], machine unlearning [5, 52] has emerged to eliminate the 031 influence of training data from a trained model [59]. The intuitive goal is to forget the specific data as if the model had never used it during training [4]. To achieve that, a direct way [52] is to retrain the 033 model from scratch by excluding the data to be unlearned, termed *exact unlearning*. Considering 034 the intensive computational cost, much attention has been paid to approximate unlearning [17, 29, 6, 11], which adjusts the trained model for approximating the behaviors of the retrained one. Prior works [29, 55] tried to conduct random-wise forgetting to unlearn uniformly sampled instances, 037 while it conflicts with the goal of model generalization and is philosophically infeasible [51]; and 038 feature-wise forgetting [58] to revoke sensitive features, but is limited to tabular data [46]. Focusing target granularity as semantic clusters, recent studies [37, 30, 6, 11] showed that *class-wise* unlearning is effective in forgetting the knowledge of a training class, either through reverse optimization [54, 29] 040 on the class data or fine-tuning on the remaining data [17] to realize catastrophic forgetting [2, 31]. 041

1053 In this work, we decouple the target concept with the class label, to model the unlearning scenarios for research explorations. To be specific, we consider the different label domains of the forgetting

⁰⁴² Despite the promising achievements, the previously studied scenario [58, 17, 6, 30, 11] mainly 043 assumed the target concept to coincide with the class label, overlooking that the practical unlearning 044 request [3, 23, 34] may violate the taxonomy of the pre-training tasks. Raised by the model users, the reported cases to be unlearned can involve different concerns from original tasks, spanning from privacy, fairness, copyright, or the hazardous capabilities [39]. The conventional matched scenario 046 is that all the identified forgetting data correspond to one pre-training class. However, those fully 047 identified cases may be only a semantic subset within a class, for which the model developer needs to 048 unlearn the small set considering reserving model utility on the other parts. Nevertheless, sometimes the user would identify limited cases of the target concept. With a conservative attitude for protecting the reputation of serving [3, 39] (e.g., IP conflicts), the developer tends to unlearn a larger semantic 051 cluster when those instances are from the same class or, more critically, across different classes. 052



Taking the *CIFAR-100* [35] dataset with its classes and superclass information, we instantiate four unlearning tasks given the same forgetting data with the class labels of "boy" and "girl": a) *all matched forgetting* (conventional scenario): unlearn "boy" and "girl" with the model trained on the classes; b) *target mismatch forgetting*: unlearn "people" with the model trained on the classes; c) *model mismatch forgetting*: unlearn "boy" and "girl" with the model trained on the superclass; d) *data mismatch forgetting*: unlearn "people" with the model trained on the superclass.

Figure 1: Illustrations of decoupling the class label and the target concept.

data \mathcal{L}_D , the model output \mathcal{L}_M , and the target concept \mathcal{L}_T in unlearning. We introduce two relations between two label domains, i.e., \mathcal{L}_1 matches \mathcal{L}_2 ($\mathcal{L}_1 = \mathcal{L}_2$) and \mathcal{L}_1 is the subclass domain of \mathcal{L}_2 ($\mathcal{L}_1 \prec \mathcal{L}_2$)¹, then identify scenarios corresponding to the target concept being larger or smaller than the class unit. Assuming that the reported forgetting data should be included in the target concept, e.g., $\mathcal{L}_D \preceq \mathcal{L}_T$, we have all matched $\mathcal{L}_D = \mathcal{L}_T = \mathcal{L}_M$; target mismatch $\mathcal{L}_D = \mathcal{L}_M \prec \mathcal{L}_T$; model mismatch $\mathcal{L}_D = \mathcal{L}_T \prec \mathcal{L}_M$; and data mismatch $\mathcal{L}_D \prec \mathcal{L}_T = \mathcal{L}_M$. We further illustrated in Figure 1 using the CIFAR-100 [35] dataset to instantiate four unlearning tasks with the classes and superclass.

079 Given the aforementioned tasks, we identify new challenges with the mismatched label domains 080 (refer to Figure 2). Unlike the accurate unlearning approximation in the conventional all matched 081 task [17, 30, 6], the representative unlearning methods [58, 54] exhibit different performance gap with the retrained reference in the other tasks. Specifically, the under-entangled feature representation (when $\mathcal{L}_M \prec \mathcal{L}_T$) or the under-representative forgetting data (when $\mathcal{L}_D \prec \mathcal{L}_T$) results in insufficient 083 forgetting, while the entangled feature representation (when $\mathcal{L}_T \preceq \mathcal{L}_M$) prevents the decomposition 084 of target concept with the retaining part. The former requires target identification in the remaining 085 dataset, while the latter requires explicit target separation over the entangled feature representation. Through exploration of forgetting dynamics (refer to Figure 3), we demonstrate the feature distance 087 reflected by representation gravity is a crucial factor for the feasibility of these unlearning tasks.

Based on the above analysis, we propose a novel framework, namely, TARget-aware Forgetting 089 (TARF), for unlearning. In general, we consider two parts (refer to Eq. 5), i.e., annealed forgetting and 090 target-aware retaining, which collaboratively enable the target identification and separation for these 091 forgetting tasks. Specifically, the algorithmic framework (refer to Figure 4) incorporates an annealed 092 gradient ascent and target-aware gradient descent in a dynamical manner, which can be viewed as three phases. The first actively unlearns the identified forgetting data, and constructs the contrast 094 information to filter out the remaining data which is hard to be affected. Then, simultaneously learning the selected retaining data with gradient descent deconstructs the entangled feature representation. 096 Ultimately, the learning objective can progressively approach standard retraining using the aligned retaining data (refer to Figure 5). We present comprehensive experiments of four unlearning tasks 098 across different datasets to demonstrate the performance of TARF, and show its application on concept forgetting with stable diffusion. The main contributions of our work can be summarized as, 099

100 101

102 103

104

068

069

- Conceptually, we introduce new settings that decouple the class label and the target concept, which investigate the label domain mismatch in class-wise unlearning (in Section 3.1).
- Empirically, we systematically reveal the challenges of restrictive unlearning with the mismatched label domains, and demonstrate that the representation gravity in forgetting dynamics is critical for achieving the forgetting target in the new tasks (in Section 3.2).
- 105 106 107

 $^{{}^{1}\}mathcal{L}_{1} \prec \mathcal{L}_{2}$: For any label $y \in \mathcal{L}_{1}$, there exist label $y' \in \mathcal{L}_{2}$ that instance being labeled with y can also being labeled with y', but not all instance being labeled with y' can be labeled with y.

- Technically, we propose a general framework, namely, *TARF*, to realize the target identification and separation in unlearning. It consists of annealed forgetting and target-aware retaining which collaboratively approximate retraining on the retaining data (in Section 3.3).
 - Experimentally, we conduct extensive explorations to validate the effectiveness of our framework and perform various ablations to characterize algorithm properties (in Section 4).

2 Preliminaries

108

110

111

112

113 114

115 116

117

118

119

132

133

134

In this section, we briefly introduce the problem settings of class-wise machine unlearning, and compare the differences between ours and the conventional setting considered in the previous research works. More details about the unlearning baselines considered in our work can refer to Appendix C.1.

120 **Problem setup.** Following the literature [52, 59], we mainly consider the multi-class classifica-121 tion [18] as the original training task for class-wise unlearning. Let $\mathcal{X} \subset \mathbb{R}^d$ denote the input space 122 and $\mathcal{Y} = \{1, \ldots, C\}$ denote the label space, where *C* is the number of classes, the training dataset 123 $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ consists of two subsets in machine unlearning, e.g., the forgetting dataset \mathcal{D}_f and 124 the retaining dataset $\mathcal{D}_r = \mathcal{D} \setminus \mathcal{D}_f$. Building upon the model $f_{\theta^*} : \mathcal{X} \to \mathcal{Y}$ trained on \mathcal{D} with the loss 125 function ℓ , the general goal of this problem is to find an unlearned model θ^*_{un} , which approximates 126 the behaviors of the model θ^r that retrained on \mathcal{D}_r from scratch,

$$\theta_{un}^{*} = \arg\min_{\theta} \frac{1}{|\mathcal{D}|} \sum_{(x,y)\sim\mathcal{D}} \mathcal{R}(\theta, \theta^{r}, x, y) \quad \text{s.t. } \theta^{r} = \arg\min_{\theta} \underbrace{\frac{1}{|\mathcal{D}_{r}|} \sum_{(x,y)\sim\mathcal{D}_{r}} \ell(f_{\theta}(x), y)}_{L_{\text{retrain}}}, \tag{1}$$

where \mathcal{R} indicates a general risk measure for model behavior consistency [17, 52], which can be instantiated by comprehensive evaluation metrics [30, 11] (e.g., unlearning accuracy (UA), retaining accuracy (RA), and others related to privacy) in experiments to pursue the unlearning efficacy and the model utility [59]. The specific definition of evaluation metrics can be referred to in Section 4.1.

Dataset partition in mismatched setting. As the target concept is decoupled from the class label, we adopt D_t to indicate the dataset of the target concept, D_f to indicate the *given forgetting* dataset, and summarize the scenarios in Table 1. We can find that the previous assumptions of $D_f = D_t$ and $D_r = D \setminus D_f$ only hold in all matched setting. In model mismatch forgetting, the former is still

Table 1: Training set data partition.

Scenario	Data Partition				
All matched Target mismatch Model mismatch Data mismatch	$ \left \begin{array}{c} \mathcal{D}_f = \mathcal{D}_t \\ \mathcal{D}_f \subset \mathcal{D}_t \\ \mathcal{D}_f = \mathcal{D}_t \\ \mathcal{D}_f \subset \mathcal{D}_t \end{array} \right. $	$ \begin{vmatrix} \mathcal{D}_{un} = \mathcal{D}_r \\ \mathcal{D}_{un} = \mathcal{D}_{fr} \cup \mathcal{D}_r \\ \mathcal{D}_{un} = \mathcal{D}_r \\ \mathcal{D}_{un} = \mathcal{D}_{fr} \cup \mathcal{D}_r \end{vmatrix} $			

held while we notice that there exists *affected retaining* data (like the left of Figure 2) in \mathcal{D}_{ar} having the same class label with that in \mathcal{D}_{f} ; in target/data mismatch forgetting, $\mathcal{D}_{f} \subseteq \mathcal{D}_{t}$ and the remaining dataset $\mathcal{D}_{un} = \mathcal{D} \setminus \mathcal{D}_{f}$ include both true retaining dataset $\mathcal{D}_{r} \subseteq \mathcal{D}_{un}$ and the *false retaining* dataset (like the right of Figure 2) $\mathcal{D}_{fr} = \mathcal{D}_{t} \setminus \mathcal{D}_{f}$, where the data belong to the target concept but have not been unidentified. Considering task feasibility, we assume that the number of classes in \mathcal{D}_{un} belonging to the target concept is known in target mismatch forgetting, and the retrained models are all trained using $\mathcal{D}_{r} = \mathcal{D} \setminus \mathcal{D}_{t}$. Details of scenario construction and class information can refer to Appendix D.4.

Different focus from prior methods. Existing studies [30, 6] generally assume that $\mathcal{D}_f = \mathcal{D}_f$ and 149 $\mathcal{D}_{\rm r} = \mathcal{D} \setminus \mathcal{D}_{\rm f}$. The common approximation unlearning methods either focus on retaining or forgetting 150 objectives. The former, represented by Fine-tuning (FT) [58], fine-tunes the model θ^o on \mathcal{D}_r to induce 151 catastrophic forgetting over \mathcal{D}_{f} . Later advances assign random labels [17] on \mathcal{D}_{f} to enforce forgetting 152 or adopt L_1 -norm [30] to infuse weight sparsity in approximation. The latter, represented by gradient 153 ascent (GA), reverse gradient updates on \mathcal{D}_{f} . And another line of works [29] utilizes the influence 154 function [33] to erase the influence. More recently, adversarial perturbation [6] on \mathcal{D}_{f} is employed to shrink the decision boundary for the target class. From a different perspective, we explore the label 156 domain mismatch that relaxes the previous assumption. More discussions are in Appendixes B and C. 157

158

3 TARF: TARget-aware Forgetting

159 160

161 In this section, we first introduce the motivation and reveal the challenges of restrictive unlearning when the class label and the target concept are decoupled (Section 3.1). Second, we present systematic



We conduct various unlearning methods for the four tasks. In conventional all-matched forgetting, all the methods can perform similarly to Retrained. In contrast, we can find that model mismatch forgetting can be affected by the trained model, coupling the behaviors on the forgetting and affected retaining (under the same superclass) data, leaving less accuracy gap between them. In target or data mismatch forgetting, the class labels can not fully represent the target concepts, leaving the non-zero accuracy on the false retaining data (belongs to the target concept).

Figure 2: The challenges of restrictive unlearning with the mismatched label domains.

exploration through the perspectives of feature representation and forgetting dynamics (Section 3.2). Lastly, we propose our novel framework, i.e., *TARget-aware Forgetting* (TARF) (Section 3.3).

3.1 MOTIVATION: EXPLORING MISMATCHED TAXONOMY IN UNLEARNING

Given its technical nature of mitigating the data influence from a trained model, unlearning is given 189 broader significance in the context of trustworthiness [3], where the requests can be varied beyond 190 the withdrawal from data owner [51], and may be applied in mitigating bias [60] to improve fairness, 191 erasing harmful content [39] to ensure safety usage, or removing inappropriate content [14] for social 192 good. Recently, a series of studies [17, 58, 30, 11, 6] have several proposals on forgetting a training 193 class of the models, and demonstrated it can be successfully achieved by partially scrubbing the class 194 data or fine-tuning on the retaining data to realize catastrophic forgetting [13, 19]. However, a general 195 scenario considered in previous works is that the target concept is aligned with the taxonomy of the 196 pre-training tasks, which may not always hold in practical scenarios with the previous meanings (we 197 leave more discussion in Appendix D.5). This naturally motivates the following research question,

What if the class labels and target concept do not coincide in unlearning?

To investigate it, we consider different label domains on \mathcal{L}_D of forgetting data, \mathcal{L}_M of model outputs, and \mathcal{L}_T of target concepts. Assuming that $\mathcal{L}_D \preceq \mathcal{L}_T$, we have either $\mathcal{L}_D = \mathcal{L}_T$ or $\mathcal{L}_D \prec \mathcal{L}_T$. When $\mathcal{L}_D = \mathcal{L}_T$, we have all matched if $\mathcal{L}_D = \mathcal{L}_M$ (e.g., forgetting "boy" and "girl" with the model trained on classes) and model mismatch if $\mathcal{L}_T \prec \mathcal{L}_M$ (e.g., forgetting "boy" and "girl" with the model trained on superclass); When $\mathcal{L}_D \prec \mathcal{L}_T$, we have target mismatch if $\mathcal{L}_D = \mathcal{L}_M$ (e.g., forgetting "people" with the model trained on classes given "boy" and "girl") and data mismatch if $\mathcal{L}_T = \mathcal{L}_M^2$ (e.g., forgetting "people" with the model trained on superclass given "boy" and "girl").

In Figure 2, we conduct the unlearning on the four forgetting tasks as instantiated in Figure 1. As a result, those unlearning methods, e.g., the representative methods FT, GA, and the recent L_1 -sparse [30] and BS [6] show different performance gaps compared with the retrained models except in the conventional all matched setting. It can be found that the *false retaining* data (which belong to the target concept but are not identified) are under-represented by the given forgetting data when $\mathcal{L}_D \prec \mathcal{L}_T$, as demonstrated by the non-zero accuracy on target concept in the right-bottom of Figure 2; and the *affected retaining* data (which is under the same superclass as the model trained

214 215

179

181

182 183

185

186 187

188

²Note that we also provide a detailed discussion in appendix D for all the potential cases, while some (e.g., $\mathcal{L}_T \prec \mathcal{L}_D$) are impractical and some (e.g., $\mathcal{L}_M \prec \mathcal{L}_T$) are similar to the major scenarios considered here.

218

219

220

221

223

225 226

227

228

229

230

231

232

233

235

236

237 238

239 240

241

242 243

244

245

247

248

249

260 261 262

0.20 20 0.1 50.10 Dist. Dist <u>б</u>о.о Concept-aligned Data emaining Data 40 21 20 22 43 16 54 -20 0 20 Pre-trained 40 -20 0 Pre-trained 40 Class Indexes (a) Distances in the model trained by classes (b) Distances in the model trained by superclass 20.0 17.5 15.0 13.0 12.5 10.0 7.5 Value acv racy Forgetting Data Forgetting D . Concept-aligned Data -Class-aligned Data Accu Loss / Accu Remaining Data Remaining Data 40 5.0 20 2.5 0.0 Epoch Epoch Epoch Epoch (c) Forgetting with under-entangled representation (d) Forgetting with entangled representation

We present the tSNE visualization of the learned features from the pre-trained model trained by (a) classes and (b) superclass, and their averaged Euclidean distance to cluster centers of the forgetting data. At the bottom, we show the averaged loss/accuracy value of forgetting data, concept/class-aligned data, and the remaining data during GA on the two representations. Note that we only show the 5 classes in tSNE due to the large number of remaining classes, we also provide the full results of the unlearned representations in Appendix F.2.

Class-alg



on) are entangled with the forgetting part when $\mathcal{L}_T \prec \mathcal{L}_M$, as demonstrated by the less accuracy gap between forgetting and affected retaining data than that of Retrained in the left-bottom of Figure 2.

SYSTEMATIC EXPLORATION ON FORGETTING DYNAMICS 3.2

The mismatch of label domains affects the construction of model representation in unlearning, which requires us to explore it further to understand the underlying mechanism of the performance gaps. In 246 Figure 3, we take a closer look at the relation between the representation and forgetting dynamics, for which we use the analogy "gravity", originating in physics, as an intuitive inspiration for describing the mutual influence on forgetting data and the other data, specializing in the following ability of unlearning behaviors with the representation similarity, similar to the mutual attraction with objects. 250

Target or data mismatch. In both target and data mismatch forgetting, we have $\mathcal{L}_D \prec \mathcal{L}_T$, which means the forgetting data is a subset of the target concept, i.e., $\mathcal{D}_f \subset \mathcal{D}_f$. As indicated in Figures 3(a) 253 and 3(c), partially relying on the forgetting or remaining data can not fully represent the target concept, 254 and leaves non-zero accuracy on false retaining data. We can summarize the following observation,

255 **Observation 3.1** (Insufficient representation). Given $\mathcal{L}_D \prec \mathcal{L}_T$ that indicates $\mathcal{D}_f \subset \mathcal{D}_t$, and a cluster 256 center h^* of feature representations $h_{x \sim \mathcal{D}_{\mathrm{f}}}(x)$ extracted at the pre-trained model θ , as well as a 257 distance measure $d(\cdot, \cdot)$, the sample $(x^u, y^u) \sim (\mathcal{D}_t \setminus \mathcal{D}_f)$ exhibits weak gravity following the sample 258 $(x,y) \sim \mathcal{D}_{f}$ on the forgetting dynamic $\epsilon = \mathbb{E}(\ell(f_{\theta}(x), y) - \ell(f_{\theta}(x), y))$ with a large value ζ_{1} , under 259 the observation interval t from θ to the unlearned model θ^t ,

$$(d(h(x^{u}), h^{*}) > \sup_{x \sim \mathcal{D}_{f}} d(h(x), h^{*})) \implies |(\ell(f_{\theta}(x^{u}), y^{u}) - \ell(f_{\theta^{t}}(x^{u}), y^{u})) - \epsilon| > \zeta_{1}.$$
(2)

Model mismatch forgetting. In this task, we have $\mathcal{L}_D = \mathcal{L}_T$ while $\mathcal{L}_T \prec \mathcal{L}_M$. Regarding the model trained by the superclass, it can be found in Figures 3(b) and 3(d) that the features of forgetting 264 data and affected retaining data are closely entangled, showing that the unlearning of the forgetting 265 data can unavoidably affect the representation of the other part. In contrast, it is also notable in the 266 left-bottom of Figure 2 that the accuracy gap between forgetting data and affected retaining data is 267 expected to be large in the retrained reference. We summarize the observation as follows, 268

Observation 3.2 (Decomposition lacking). Given $\mathcal{L}_T \prec \mathcal{L}_M$ that indicates the broader representation 269 region for $\mathcal{D}_z := \mathcal{D}_{ar} \cup \mathcal{D}_t$ within the same class z, where the sample $x^m \sim (\mathcal{D}_z \setminus \mathcal{D}_t)$ exhibits strong



The overall framework consists of two objective parts, e.g., annealed forgetting and target-aware retaining, which can be regarded as three phases to enable all the class-wise unlearning tasks through the view of the unlearning process. (a) Phase I utilizes the gradient ascent to construct dynamic information for all class data; (b) Phase II simultaneously considers gradient ascent on forgetting data and gradient decent on remaining data that is hard to affect to separate target concept; (c) Phase III conducts gradient descent on the selected data to approximate the retraining.

Figure 4: Overview of the proposed framework TARF.

gravity following the forgetting dynamic $\epsilon = \mathbb{E}(\ell(f_{\theta}(x), y) - \ell(f_{\theta}(x), y))$ with a small value ζ_2 ,

$$(d(h(x^m), h^*) \le \sup_{x \sim \mathcal{D}_z} d(h(x), h^*)) \implies |(\ell(f_\theta(x^m), y^m) - \ell(f_{\theta^t}(x^m), y^m)) - \epsilon| < \zeta_2.$$
(3)

Forgetting dynamics with representation distance. Despite the issues revealed by previous observations under label domain mismatch, the forgetting performance varying obviously on different representations also provides clues on addressing them. Notably, we can find that GA achieves better forgetting efficacy on the data mismatch forgetting as the feature representation of the forgetting data and false retaining data is entangled. Through the effect of actively forgetting the given data on the other parts of data, we have the following that reveals a crucial factor for achieving these tasks,

300 **Definition 3.3** (Representation gravity). Given the empirically demonstrated Observation 3.1 and Ob-301 servation 3.2, we can have an indicator $I_{con}(x, y, \theta)$ to reflect the representation similarity $d(h(x), h^*)$ in the model θ^t , which is a crucial factor for the feasibility of unlearning under mismatch, 302

$$I_{\rm con}(x,y,\theta) = |\ell(f_{\theta}(x),y) - \ell(f_{\theta^t}(x),y)|.$$

$$\tag{4}$$

It is empirically demonstrated in Figure 3 where we present the average representation distance 305 to the cluster center of forgetting data, the corresponding changes in accuracy and loss values 306 show that the smaller the distance in representation level, the similar forgetting dynamics the 307 model would have on prediction. As the Observations 3.1 and 3.2 reveal the issues of insufficient 308 representation and decomposition missing, we can utilize the representation gravity to identify the 309 other unidentified forgetting data in the remaining set, and reveal the needs of deconstructing the 310 pre-entangled representation by simultaneously considering the forgetting and retaining objectives. 311

312 3.3 ALGORITHM FRAMEWORK OF TARF 313

286

287

288

289 290

291 292

293

295

296

297

298

299

303 304

317 318

319 320

314 Based on the previous analysis, we introduce the whole framework of *TARget-aware Forgetting* 315 (TARF), to enable the four class-wise unlearning tasks. Given the identified forgetting data, we illustrate the overall process in Figure 4, and introduce its dynamic learning objective as follows: 316

$$L_{\text{TARF}} = \underbrace{k(t) \cdot \left(-\frac{1}{|\mathcal{D}_{f}|} \sum_{(x,y) \sim \mathcal{D}_{f}} \ell(f(x), y)\right)}_{\text{Annealed Forgetting } L_{f}(k)} + \underbrace{\frac{1}{|\mathcal{D}_{\text{un}}|} \sum_{(x,y) \sim \mathcal{D}_{\text{un}}} \ell(f(x), y) \cdot \tau(x, y, t),}_{\text{Target-aware Retaining } L_{u}(\tau)}$$
(5)

321 where k(t) serves as an annealing strategy to control the strength of the forgetting part. Along with 322 training, we expect the overall objective to approximate the retraining ones $L_{\text{TARF}} \rightarrow L_{\text{retrain}}$ through, 323 L

$$f(k) \to 0, \quad L_{u}(\tau) \to L_{retrain},$$
(6)



We show (a) accuracy changes in target identification in target mismatch forgetting, unlearning performance using different forgetting classes in data mismatch forgetting; (b) accuracy gap of retaining and forgetting part of the same class, as well as the need of reconstruction.

Figure 5: Target identification and target separation for unlearning under mismatch.

given the initially provided forgetting data \mathcal{D}_{f} and the remaining set \mathcal{D}_{un} . Specifically, we design the two dynamic hyperparameters k(t) and $\tau(x, y, t)$ as follows to achieve that,

$$k(t) = \frac{k \cdot (T - t - t_0)}{T}, \quad t \in [0, T]; \quad \tau(x, y, t) = \begin{cases} 0 & I_{\text{con}}(x, y, \theta_{t_1}) > \beta \text{ or } t < t_1, \\ 1 & I_{\text{con}}(x, y, \theta_{t_1}) < \beta \text{ and } t \ge t_1, \end{cases}$$
(7)

where T indicates the total training time (e.g., epochs), and the value of k(t) decreases with the training process, β can be estimated by the information of forgetting dynamics about the specific unlearning request and the rank of loss/accuracy change at t_1 , t_0 and t_1 respectively control the end time of active forgetting and the begin time of retaining part. The overall process can be regarded as,

Phase I: Target Identification. Before t_1 , since $\tau(x, y, t) = 0$, Eq. 5 can be formalized as, $L_{\text{TARF-Phase-I}} = k(t) \cdot \left(-\frac{1}{|\mathcal{D}_t|} \sum_{(x,y)\sim\mathcal{D}_f} \ell(f(x), y)\right)$, in which the retaining part is waiting for the dynamic information revealed by this phase. As shown in Figures 3, the false retaining data in \mathcal{D}_{fr} can be identified due to the similar forgetting dynamics with the forgetting data. To conduct the specific selection, we utilize the class label information in our main tasks, for which we detail the implementation of controlling β in Appendix E. In Figure 5(a), we show the selected classes via the accuracy variance and validate identification efficacy with varied given forgetting classes.

Phase II: Target Separation. After phase I, the retaining part is engaged with the forgetting part with the identified data \mathcal{D}_{fr} and the remaining retaining data \mathcal{D}_{r} . By simultaneously considering the forgetting and retaining part as Eq. 5, $L_{TARF-Phase-II}$ encourages the model to deconstruct the target concept and reconstruct the feature representation of the retaining part, which can effectively decouple the pre-entangled feature in the model mismatch forgetting. In the left of Figure 5(b), we compare the accuracy gap on the retaining and forgetting part, demonstrating the necessity of considering two parts of objectives to separate the entangled feature representation in model mismatch forgetting.

Phase III: Retraining Approximation. After t_0 , we focus on retaining in the current phase, which approximates the retraining objective as follows, $L_{\text{TARF-Phase-III}} = \frac{1}{|\mathcal{D}_{un}|} \sum_{(x,y) \sim \mathcal{D}_{un}} \ell(f(x), y) \cdot \tau(x, y, t)$, where we use τ at t_1 to indicate the identified hard-to-effect retaining data, and continually reconstruct the representations. Since the general goal of unlearning considered in our work is similar to retraining, this phase can prevent excessive forgetting. In the right of Figure 5(b), we compare the performance using different lengths of this phase for approximating the retrained reference.

371 372

334 335

336 337

338

344

345

346

347 348

349

350

351

352

353

354

355

4 EXPERIMENTS

373

In this section, we present the comprehensive experimental results. To begin with, we introduce the
basic setups for the unlearning and the evaluation (Section 4.1). Then, we validate the effectiveness
of our method in four unlearning scenarios with the decoupled class label and the target concept
(Section 4.2). To better understand its properties, we conduct various experiments on the ablation study
and provide further discussions (Section 4.3). More details and results are provided in Appendix F.

379

380

381

408

409

Type / \mathcal{D}	Dataset			CIF	AR-10					CIF	AR-100		
	Method / Metrics	UA	RA	TA	MIA	Gap↓	TIME↓	UA	RA	TA	MIA	Gap↓	TIME↓
	Retrained (Ref.)	0.00	99.51 98.62	94.69 92.36	100.00	-	43.3	0.00	97.85	76.03	100.00	-	43.2
	RL [56]	4.13	97.65	91.23	100.00	2.36	4.88	1.00	96.09	72.00	100.00	1.70	4.96
All matched	GA [28] IU [29]	0.49 0.22	95.24 88.15	88.17 82.38	99.78 99.96	2.88 5.99	0.25 0.45	0.00	94.74 37.61	68.56 29.58	99.89 100.00	3.01 26.67	0.06 0.51
An matcheu	BS [6]	25.04	87.94 94 20	80.90	88.67 100.00	15.43	0.82	4.60	90.18 94.60	63.66 71.57	99.55 100.00	6.27	0.78
	SalUn [11]	0.00	91.32	86.87	100.00	4.00	5.65	0.00	75.34	62.14	100.00	9.10	5.75
	SCRUB [37]	0.00	99.94	91.00	100.00	1.03	2.88	0.00	99.98	76.75	100.00	0.71	3.23
	IARF (ours)	0.00	98.23	91.95	100.00	1.01	4.21	0.00	96.90	72.53	100.00	1.11	4.68
	FT [58]	87.76	99.58 98.53	95.91 93.56	20.57 9.56	5.33	43.8 4.29	88.22	98.58 95.02	78.50	16.33	4.58	43.8 4.86
	RL [56] GA [28]	53.69 5.76	97.85 86.99	92.39 82.20	96.60 94.98	$28.84 \\ 45.68$	4.82 0.25	80.11	95.83 94.83	79.83 76.96	99.00 97.78	21.35 39.68	4.93 0.06
Model	IU [29]	23.69	87.34	82.57	89.87	39.74	0.44	34.67	96.83	79.08	86.44	29.14	0.49
mismatch	L_1 -sparse [30]	93.11	50.77 94.76	49.39 91.63	95.96 14.44	62.05 5.15	0.79 4.24	90.22	95.90 94.78	72.28	95.22 18.88	37.14 3.25	0.89 5.00
	SalUn [11]	8.91 95.14	93.95 99.81	84.38 94 22	99.32 15.38	43.69 3.61	6.04 3.06	66.33	78.83 99.74	70.78 79.23	77.00 21.11	25.15 2.45	5.97 4.12
	TARF (ours)	91.11	97.49	92.49	17.82	2.90	4.31	86.67	97.05	80.07	26.00	1.21	4.81
	Retrained (Ref.)	0.00	99.38	93.85	100.00	-	52.1	0.00	97.85	73.72	100.00	-	53.2
	FT [58] RI [56]	50.43	98.47 97.56	91.65 90.90	50.44 56.23	25.78 24.95	4.38	58.18	96.32 96.05	72.53	46.76 46.98	28.54 28.81	5.00
Target	GA [28]	40.82	97.01	89.51	64.32	20.80	0.26	21.38	96.64	70.22	90.67	8.86	0.05
mismatch	BS [6]	53.62	88.65	75.39	58.75 76.33	26.62	0.44	40.44	98.32	29.38 68.66	85.16	42.93	0.30
	L_1 -sparse [30] SalUn [11]	49.47	93.61 91.08	88.83 86 31	51.24 60.94	27.26 25.38	4.38 5.90	56.09 59.64	94.63 75.52	72.00 62.37	48.04 65.96	28.25 27.35	4.78 5.81
	SCRUB [37]	49.98	99.94	92.10	50.18	25.53	2.89	59.64	99.99	75.32	44.89	29.90	3.52
	TARF (ours)	0.06	97.57	90.81	100.00	1.23	4.23	0.31	97.35	73.68	100.00	0.21	4.85
	Retrained (Ref.)	0.00	99.54 08.40	95.56	100.00	-	52.1	0.00	98.50 95.66	80.15	100.00	37 15	53.2
	RL [56]	76.47	97.68	91.93	49.81	33.04	4.76	89.78	96.82	79.90	70.76	30.49	4.97
Data	GA [28] IU [29]	8.69 22.84	96.41 95.50	90.78 89.54	93.03 88.57	5.89 11.08	0.25 0.44	6.00	97.65 98.96	79.23 78.20	98.04 88.09	2.43 11.46	0.05 0.48
mismatch	BS [6]	16.70	61.21	49.76	92.24	22.37	0.82	15.38	98.50	72.28	96.22	6.76	0.96
	SalUn [11]	51.77	93.87	90.46	63.52	24.75	5.72	72.93	78.87	71.04	54.13	36.89	5.72
	SCRUB [37]	97.13	99.89	95.03	10.99	46.76	2.94	95.50	99.79	79.68	15.11	45.54	3.68
	TARF (ours)	0.00	98.17	93.09	100.00	0.96	4.22	0.00	95.01	78.98	100.00	1.17	4.78

Table 2: Main Results (%). Comparison with the unlearning baselines. All methods are trained on the same backbone, i.e., the basis of unlearning initialization is the same (except for the reference Retrained from scratch). Bold numbers are superior results. ↓ indicates smaller values are better. (The complete results under multiple runs are summarized with mean and std values in Appendix F.7)

4.1 EXPERIMENTAL SETUP

410 **Datasets, models, and baselines.** In our experiments, we explore machine unlearning for conven-411 tional image classification tasks. Since the introduced unlearning settings need a coarse-to-fine label 412 structure, we adopt the benchmarked dataset, e.g., CIFAR-10/CIFAR-100 [35] with their superclass 413 information in the major experiments. Specifically, we train two models based on the original classes 414 and its superclass respectively, and instantiate four tasks (as illustrated in Figure 1) of unlearning 415 with the decoupled class label and the target concept. The detailed information is summarized in Appendix D.4. Following previous works [58, 30, 11], we use ResNet-18 [24] as the major architecture 416 to obtain the original trained models with standard learning [18], and then set it to be the basis for 417 unlearning. In addition, we also adopt TinyImageNet [38], ImageNet [36] for large-scale experiments. 418 As for comparison, we consider four representative baselines with the retrained model (Retrained), 419 e.g., FT [58, 18], RL [56], GA [28], IU [29], and also consider four recent advanced methods, e..g., 420 BS [6], L₁-sparse [30], SalUn [11], and SCRUB [37]. The method details are in Appendix C.1. 421

422 **Evaluation metrics.** The general target of class-wise unlearning considered in this work is to 423 approximate the Retrained model [58]. To give a comprehensive evaluation, we adopt 5 specific 424 evaluation metrics in classification tasks following previous works [30, 11]. We utilize Unlearning 425 Accuracy (UA) to evaluate the accuracy of the unlearning targeted subset; Retaining Accuracy (RA) 426 to evaluate the accuracy of the retaining subset; Testing Accuracy (TA) to evaluate the generalization 427 ability of the model; Membership Inference Attack (MIA) to evaluate the efficacy of unlearning by 428 the confidence-based predictor. Note that any single indicator does not represent optimally in the approximation of a Retrained reference. All the above will be compared with that of the Retrained 429 model and summarized in a "Gap" value (averaged gap with Retrained) to indicate the overall 430 performance (the lower the better), and we also adopt TIME to present the computational time. The 431 implementation of evaluation metrics in different unlearning scenarios is detailed in Appendix C.2.

Type / \mathcal{D}	Dataset			All r	natched					Model	mismatch	I	
	Method / Metrics	UA	RA	TA	MIA	Gap↓	TIME↓	UA	RA	TA	MIA	Gap↓	TIME↓
	Retrained (Ref.) FT [58] RL [56] GA [28] L ₁ -sparse [30] SCRUB [37]	0.00 3.80 73.20 5.70 3.70 0.00	74.32 77.66 69.87 63.26 76.63 75.06	63.13 62.98 60.49 57.09 62.55 63.82	100.00 97.30 18.40 87.50 97.50 100.00	2.50 40.47 8.83 2.28 0.36	217.0 30.41 225.13 0.34 40.79 66.69	34.80 59.30 84.10 6.30 59.40 37.70	71.26 77.26 68.53 63.17 76.30 73.89	64.29 62.92 60.63 58.04 62.80 64.20	66.90 38.00 8.00 90.70 38.80 57.30	- 15.19 28.64 16.66 14.81 3.81	256.14 37.44 226.79 0.34 37.05 58.53 28.21
Tiny ImageNet	Dataset Method / Metrics	UA	RA	Target	matched MIA	Gap↓	TIME↓	UA	RA	Data 1 TA	nismatch MIA	Gap↓	TIME
	Retrained (Ref.) FT [58] RL [56] GA [28] L ₁ -sparse [30] SCRUB [37]	0.00 29.67 68.93 11.33 28.93 25.67	72.83 75.94 69.97 63.63 75.18 75.31	65.12 62.97 60.55 57.26 62.55 63.85	100.00 69.30 22.00 81.00 69.60 73.80	- 16.41 38.59 11.85 16.06 13.90	213.05 30.41 225.13 0.34 40.79 66.69	0.00 64.33 84.27 7.33 63.90 44.07	71.37 75.45 68.64 63.44 74.80 74.02	65.76 62.96 60.59 58.24 62.80 64.25	100.00 30.60 7.86 89.80 31.30 46.93	35.15 46.08 8.25 34.75 25.33	252.62 37.44 226.79 0.34 37.05 58.53
	TARF (ours)	5.07	75.78	62.72	97.53	3.22	28.81	0.00	74.85	62.59	100.00	1.66	27.92

Table 3: Results (%). Comparison with the baselines on TinyImageNet-200 trained on a larger model
 structure, i.e., ResNet101. (More results on large-scale dataset can refer to Appendix F.5)

4.2 PERFORMANCE EVALUATION

In this part, we present the main comparison results with those considered baselines in the four unlearning tasks, evaluated with the four detailed metrics and the overall performance gap with retrained references. We also report results under multiple runs with mean and std values in Appendix F.7.

As the performance reference, all the retrained models 453 (termed Retrained) are trained with the fully aligned re-454 taining data. Note that the UA of Retrained (Ref.) in 455 the model mismatch scenario is not equal to 0 since it is 456 evaluated with superclass label. In Table 2, we can find 457 the previous unlearning methods achieved satisfactory per-458 formance on the conventional all matched forgetting, but 459 did not perform well on the other three newly considered 460 tasks with the label domain mismatch. Specifically, since 461 the previous methods partially rely on forgetting data or 462 remaining data, it results in ineffective or excessive forget-463 ting due to the insufficient representation or decomposition missing. For example, FT can retain a similar RA with 464 the Retrained but be less effective in forgetting, while GA 465 reaches the lowest UA across different tasks but sacrifices 466 too much model performance on the retaining dataset. In 467 contrast, our TARF can consistently perform better (or 468 comparable with the best method) through simultaneous 469 gradient ascent and target-aware gradient descent to re-470 strict the forgetting regions on the four unlearning tasks.



Figure 6: Image generation results of original and unlearned stable diffusion. More results are in Tables 15 and 16.

To verify TARF in a large-scale dataset, we also conduct 472 unlearning on Tiny-ImageNet trained on ResNet-101 in Table 3, and also in ImageNet-1k as well 473 as forgetting multiple classes in Appendix F. The results again show that our TARF can achieve 474 satisfactory performance regarding the overall gap with Retrained. In addition, we also conduct a 475 case study to unlearn the concept in stable diffusion [25], considering the practical data mismatch 476 forgetting where users report some undesirable examples to represent the unwanted concept. In 477 Figure 6, we show the efficacy of unlearning the "springer" and "tench", compared with the original 478 generation results and certain label (CL) mismatching, TARF can identify the potential samples with 479 similar features on the target concept and perform more thorough forgetting, e.g., generating less 480 feature related to "dog" and "fish" in Figure 6. The full results are provided in Appendix E.1.

481 482

483 484

471

4.3 Ablations and Further Exploration

In this part, we provide further exploration of the three class-wise unlearning tasks and conduct various ablation studies to characterize TARF. More results and discussion are provided in Appendix F.

450

451



Figure 7: Ablation studies: *Left:* performance using different initialized k on all matched forgetting; *middle-left:* effects of constant or different dynamic gradient ascent controlled by k(t); *middle*right: comparison of forgetting with different model structures; right: comparison of using different operations on the selected forgetting data. More experimental details can refer to Appendix F.

500 Weighted control on annealed gradient ascent. To analyze the annealed gradient ascent, we present the results on the left of Figure 7 to show the effects of initialized strength k on the all matched 502 forgetting task using the CIFAR-100 dataset. The results show that an appropriate k (e.g., about 0.05) can help the model to achieve a satisfactory performance. However, the larger k results in lower retaining performance and higher Gap value as the strength increases on the feature deconstruction. 504 505

506 **Constant or dynamic gradient ascent for forgetting.** In the middle-left of Figure 7, we study 507 whether we need the learning-rate-reduced k for the forgetting part. Specifically, we compare it 508 with using constant k and learning-rate-increased k on two model mismatch forgetting tasks. The 509 results demonstrate that annealed gradient ascent can achieve more similar performance with the Retrained on forgetting data. The gradient ascent is considered simultaneously with gradient descent 510 for restricting the forgetting region, while we adopt the annealed one since the unlearning target is to 511 approximate the retrained model instead of continually maximizing the loss of forgetting data. 512

- 513 514 **Unlearning on models trained by different structures.** In the middle-right of Figure 7, we 515 investigate forgetting on the models pre-trained using different structures, e.g., ResNet-18 [24], VGG-16bn [53], and WideResNet-50 [63]. The results of TARF on the model mismatch forgetting 516 demonstrate that our TARF can achieve the lower performance gap than FT, evaluated with the 517 retrained reference. With the increasing model capacity on the original training tasks, we can also find 518 the model with a smaller capacity makes it harder to decompose the entangled feature representation 519 for achieving the unlearning target, which increases the representation complexity.
 - 520 521 522

523

524

525

526

495

496

497

498 499

501

Different operations on the selected forgetting data. In the right of Figure 7, we present the ablation on the specific gradient operation on the selected forgetting data \mathcal{D}_{fr} . We compare using the gradient ascent (-k(t)) and cleaning (0) with the Retrained reference in target mismatch forgetting. Except for the similar forgetting efficacy achieved by the three trials, major differences exist in the performance evaluated by RA. The results show that gradient cleaning may be a better choice for \mathcal{D}_{fr} to not deconstruct the feature representation too much and affect the retaining accuracy.

527 528 529

530

5 CONCLUSION

531 In this work, we decouple the class label and target concept in class-wise unlearning. By introducing 532 the label domain mismatch among forgetting data, model output, and target concept, we uncover three additional tasks beyond the conventional all matched forgetting, e.g., target mismatch, model 534 mismatch, and data mismatch forgetting. We identify the insufficient representation and decomposition lacking of restrictively forgetting the target concept, and reveal the crucial forgetting dynamics 536 in the representation level for the feasibility of these unlearning requests. Based on that, we propose the TARF that assigns an annealed gradient ascent on the identified forgetting data and the normal gradient descent on the selected retaining data. By collaboratively considering the forgetting/retaining 538 target, TARF is more accurate in unlearning while maintaining the rest. We hope our work can provide new insights and draw more attention toward the practical scenarios of machine unlearning.

540 ETHICS STATEMENT

This paper does not raise any ethical concerns. This study does not involve any human subjects, practices to data set releases, potentially harmful insights, methodologies and applications, potential conflicts of interest and sponsorship, discrimination/bias/fairness concerns, privacy and security issues, legal compliance, and research integrity issues.

Reproducibility Statement

We provide the link to our source codes to ensure the reproducibility of our experimental results: https://anonymous.4open.science/r/TARF-83B5/. Below we summarize some critical aspects to facilitate reproducible results:

- **Datasets.** The datasets we used are all publicly accessible, which is introduced in Section 4.1. For our newly introduced unlearning scenarios, we provide the specific dataset construction in our code, implemented as described in Section 4.1 and Appendix D.4.
- Assumption. Following the previous work [58, 30, 11], We set our experiments to a tuning scenario where a well-trained model is available, and all the training samples are available but limited samples are labeled as "to be unlearned".
- **Open source.** The code repository will be available in an anonymous repository for the reviewing purposes. We provide a series of unlearning methods considered in our work and also the pre-trained model for unlearning.
- Environment. All experiments are conducted with multiple runs on NVIDIA Tesla V100-SXM2-32GB GPUs with Python 3.8 and PyTorch 1.8. More detailed requirements can also refer to the environment descriptions in our aforementioned source codes.

References

- [1] Jessa Bekker and Jesse Davis. Learning from positive and unlabeled data: A survey. *Machine Learning*, 109:719–760, 2020.
- [2] Dimitri P Bertsekas. Nonlinear programming. In *booktitle of the Operational Research Society*, 1997.
- [3] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. In *arXiv*, 2021.
 - [4] Lucas Bourtoule, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In 2021 IEEE Symposium on Security and Privacy (SP), pages 141–159. IEEE, 2021.
- [5] Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In 2015 IEEE symposium on security and privacy, pages 463–480. IEEE, 2015.
- [6] Min Chen, Weizhuo Gao, Gaoyang Liu, Kai Peng, and Chen Wang. Boundary unlearning. *arXiv preprint arXiv:2303.11570*, 2023.
- [7] Xuxi Chen, Wuyang Chen, Tianlong Chen, Ye Yuan, Chen Gong, Kewei Chen, and Zhangyang Wang. Self-pu: Self boosted and calibrated positive-unlabeled training. In *ICML*, 2020.
- [8] Jun Du and Zhihua Cai. Modelling class noise with symmetric and asymmetric distributions. In *AAAI*, 2015.
- [9] Marthinus du Plessis, Gang Niu, and Masashi Sugiyama. Convex formulation for learning from positive and unlabeled data. In *International conference on machine learning*, pages 1386–1394. PMLR, 2015.
- [10] Marthinus C du Plessis, Gang Niu, and Masashi Sugiyama. Analysis of learning from positive and unlabeled data. *Advances in neural information processing systems*, 27, 2014.

- [11] Chongyu Fan, Jiancheng Liu, Yihua Zhang, Dennis Wei, Eric Wong, and Sijia Liu. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. *arXiv preprint arXiv:2310.12508*, 2023.
- 597
 598 [12] Chongyu Fan, Jiancheng Liu, Alfred Hero, and Sijia Liu. Challenging forgets: Unveiling the
 599 worst-case forget sets in machine unlearning. *arXiv preprint arXiv:2403.07362*, 2024.
- [13] Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135, 1999.
 - [14] Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. *arXiv preprint arXiv:2303.07345*, 2023.
 - [15] Yingqiang Ge, Shuchang Liu, Zuohui Fu, Juntao Tan, Zelong Li, Shuyuan Xu, Yunqi Li, Yikun Xian, and Yongfeng Zhang. A survey on trustworthy recommender systems. *ACM Transactions on Recommender Systems*, 2022.
 - [16] Antonio Ginart, Melody Guan, Gregory Valiant, and James Y Zou. Making ai forget you: Data deletion in machine learning. *Advances in neural information processing systems*, 32, 2019.
- [17] Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9304–9312, 2020.
 - [18] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*. MIT Press, 2016.
 - [19] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015.
 - [20] Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac machine learning. In *Proceedings* of the AAAI Conference on Artificial Intelligence, pages 11516–11524, 2021.
 - [21] Sheng Guo, Weilin Huang, Haozhi Zhang, Chenfan Zhuang, Dengke Dong, Matthew R Scott, and Dinglong Huang. Curriculumnet: Weakly supervised learning from large-scale web images. In ECCV, 2018.
 - [22] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *NeurIPS*, 2018.
 - [23] Tatsunori Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. Fairness without demographics in repeated loss minimization. In *ICML*, 2018.
 - [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
 - [25] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
 - [26] Chris Jay Hoofnagle, Bart Van Der Sloot, and Frederik Zuiderveen Borgesius. The european union general data protection regulation: what it is and what it means. *Information & Communications Technology Law*, 28(1):65–98, 2019.
 - [27] Yu-Guan Hsieh, Gang Niu, and Masashi Sugiyama. Classification from positive, unlabeled and biased negative data. In *International Conference on Machine Learning*, pages 2820–2829. PMLR, 2019.
 - [28] Takashi Ishida, Ikko Yamane, Tomoya Sakai, Gang Niu, and Masashi Sugiyama. Do we need zero training loss after achieving zero training error? In *ICML*, 2020.
- [29] 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.

656

657

658

659

660 661

662

663

664 665

666

667 668

669

670

671

672 673

674

675

676

677

678 679

680

681

682

683 684

685

686 687

688 689

690

691

692 693

694

696

697

698 699

700

- [30] Jinghan Jia, Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, Pranay Sharma, and Sijia Liu. Model sparsity can simplify machine unlearning. In *Annual Conference on Neural Information Processing Systems*, 2023.
- [31] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 2017.
 - [32] Ryuichi Kiryo, Gang Niu, Marthinus C du Plessis, and Masashi Sugiyama. Positive-unlabeled learning with non-negative risk estimator. In *NeurIPS*, 2017.
 - [33] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International conference on machine learning*, pages 1885–1894. PMLR, 2017.
 - [34] Adriana Kovashka, Olga Russakovsky, Li Fei-Fei, and Kristen Grauman. Crowdsourcing in computer vision. In *Found. Trends Comput. Graph. Vis.*, 2016.
 - [35] Alex Krizhevsky. Learning multiple layers of features from tiny images. In *arXiv*, 2009.
 - [36] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012.
 - [37] Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. Towards unbounded machine unlearning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=OveBaTtUAT.
 - [38] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. In CS 231N, 2015.
 - [39] Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D Li, Ann-Kathrin Dombrowski, Shashwat Goel, Gabriel Mukobi, et al. The wmdp benchmark: Measuring and reducing malicious use with unlearning. In *ICML*, 2024.
 - [40] Bing Liu, Wee Sun Lee, Philip S Yu, and Xiaoli Li. Partially supervised classification of text documents. In *ICML*, pages 387–394. Sydney, NSW, 2002.
 - [41] Tongliang Liu and Dacheng Tao. Classification with noisy labels by importance reweighting. In *TPAMI*, 2015.
 - [42] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. In *booktitle of machine learning research*, 2008.
 - [43] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.
 - [44] Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C Lipton, and J Zico Kolter. Tofu: A task of fictitious unlearning for llms. arXiv preprint arXiv:2401.06121, 2024.
 - [45] Aditya Menon, Brendan Van Rooyen, Cheng Soon Ong, and Bob Williamson. Learning from corrupted binary labels via class-probability estimation. In *International conference on machine learning*, pages 125–134. PMLR, 2015.
 - [46] Vangelis Metsis, Ion Androutsopoulos, and Georgios Paliouras. Spam filtering with naive bayes-which naive bayes? In *CEAS*, volume 17, pages 28–69. Mountain View, CA, 2006.
 - [47] Seth Neel, Aaron Roth, and Saeed Sharifi-Malvajerdi. Descent-to-delete: Gradient-based methods for machine unlearning. In *Algorithmic Learning Theory*, pages 931–962. PMLR, 2021.
 - [48] Stuart L Pardau. The california consumer privacy act: Towards a european-style privacy regime in the united states. J. Tech. L. & Pol'y, 23:68, 2018.
 - [49] Jeffrey Rosen. The right to be forgotten. Stan. L. Rev. Online, 64:88, 2011.

- 702 [50] Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. Remember 703 what you want to forget: Algorithms for machine unlearning. Advances in Neural Information 704 Processing Systems, 34:18075–18086, 2021. 705 [51] Vedant Shah, Frederik Träuble, Ashish Malik, Hugo Larochelle, Michael Mozer, Sanjeev Arora, 706 Yoshua Bengio, and Anirudh Goyal. Unlearning via sparse representations. arXiv preprint 707 arXiv:2311.15268, 2023. 708 709 [52] Thanveer Shaik, Xiaohui Tao, Haoran Xie, Lin Li, Xiaofeng Zhu, and Qing Li. Exploring the 710 landscape of machine unlearning: A survey and taxonomy. arXiv preprint arXiv:2305.06360, 711 2023. 712 [53] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale 713 image recognition. In ICLR, 2015. 714
 - [54] Anvith Thudi, Gabriel Deza, Varun Chandrasekaran, and Nicolas Papernot. Unrolling sgd: Understanding factors influencing machine unlearning. In 2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P), pages 303–319. IEEE, 2022.
 - [55] Anvith Thudi, Hengrui Jia, Ilia Shumailov, and Nicolas Papernot. On the necessity of auditable algorithmic definitions for machine unlearning. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 4007–4022, 2022.
 - [56] Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. In *International Conference on Learning Representations*, 2018.
 - [57] Enayat Ullah, Tung Mai, Anup Rao, Ryan A Rossi, and Raman Arora. Machine unlearning via algorithmic stability. In *Conference on Learning Theory*, pages 4126–4142. PMLR, 2021.
 - [58] Alexander Warnecke, Lukas Pirch, Christian Wressnegger, and Konrad Rieck. Machine unlearning of features and labels. *NDSS*, 2023.
 - [59] Heng Xu, Tianqing Zhu, Lefeng Zhang, Wanlei Zhou, and Philip S Yu. Machine unlearning: A survey. *ACM Computing Surveys*, 56(1):1–36, 2023.
 - [60] Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. *arXiv preprint arXiv:2310.10683*, 2023.
 - [61] Youngsik Yoon, Jinhwan Nam, Hyojeong Yun, Jaeho Lee, Dongwoo Kim, and Jungseul Ok. Few-shot unlearning by model inversion. *arXiv preprint arXiv:2205.15567*, 2022.
 - [62] Hwanjo Yu, Jiawei Han, and KC-C Chang. Pebl: Web page classification without negative examples. *IEEE Transactions on Knowledge and Data Engineering*, 16(1):70–81, 2004.
 - [63] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *BMVC*, 2016.
 - [64] Eric Zhang, Kai Wang, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Forget-me-not: Learning to forget in text-to-image diffusion models. *arXiv preprint arXiv:2211.08332*, 2023.

715

716

717 718

719

720

721

722

723

724 725

726

727

728

729 730

731

732

733

734 735

736

737

738

739 740

741

742

A	PPEN	NDIX	
A	Rep	roducibility Statement	16
B	Disc	cussion about Related Work	16
	B .1	Machine Unlearning	16
	B.2	Positive-unlabeled Learning	16
С	Deta	ails about Considered Baselines and Metrics	17
	C .1	Unlearning Methods	17
	C .2	Evaluation Metrics Regarding Different Scenarios	19
D	Full	Discussion about Label Domain Mismatch	19
	D.1	A two-layer Label Structure of Mismatch	21
	D.2	A three-layer Label Structure of Mismatch	22
	D.3	Further Exploration on the Other 6 Different Scenarios	22
	D.4	Specific Information of the Instantiated Tasks	24
	D.5	Discussion on the Practicality of Label Domain Mismatch	25
	D.6	Discussion on the Scenario Commonalities and Framework Principles	26
£	Algo	orithm Implementation and Explanation	27
	E. 1	Case study for Unlearning Generation Concept	27
	E.2	Discussion on the Algorithm Computation Cost	28
	E.3	Discussion on TARF with Limited Class Information	29
F	Add	litional Experimental Results	32
	F.1	Extra Experimental Setups	32
	F.2	Discussion about Target Identification in unlearning	32
	F.3	Discussion about TARF on All Matched scenario	32
	F.4	Discussion about TARF on Weakly-supervised Scenario	33
	F.5	Forgetting in the Large-scale Dataset	34
	F.6	Forgetting with Different Model Structures	35
	F.7	Full Results with Different Forgetting Tasks	35
G	Bro	ader Impact	35

810 The whole Appendix is structured in the following manner. In Appendix A, we provide the anonymous 811 link to our source code and introduce the critical aspects of reproducibility. In Appendix B, we provide 812 a detailed discussion with related works of machine unlearning and other aspects. In Appendix C, 813 we review the representative baseline methods in machine unlearning, which are considered in our 814 experimental comparisons. In Appendix D, we introduce the complete scenarios considering the mismatch issues in machine unlearning, going beyond the four basic scenarios presented in the main 815 text. In Appendix E, we formally present the algorithm implementation of our proposed TARF with 816 its variant, and further explanation of the rationality of TARF in unlearning. In Appendix F, we 817 provide additional experimental results to characterize forgetting dynamics and the properties of 818 TARF. In Appendix G, we discuss the potential broader impact and limitations of our work. 819

820 821

822 823

824

825

827

828

829

830

831

832

833

834

835 836

837

A REPRODUCIBILITY STATEMENT

We provide the link to our source codes to enhance the reproducibility of our experimental results: https://anonymous.4open.science/r/TARF-83B5/. Below we summarize some critical aspects to facilitate reproducible results:

- **Datasets.** The datasets we used are all publicly accessible, which is introduced in Section 4.1. For our newly introduced unlearning scenarios, we provide the specific dataset construction in our code, implemented as described in Section 4.1 and Appendix D.4.
- Assumption. Following the previous work [58, 30, 11], We set our experiments to a tuning scenario where a well-trained model is available, and all the training samples are available but limited samples are labeled as "to be unlearned".
- **Open source.** The code repository will be available in an anonymous repository for the reviewing purposes. We provide a series of unlearning methods considered in our work and also the pre-trained model for unlearning.
- Environment. All experiments are conducted with multiple runs on NVIDIA Tesla V100-SXM2-32GB GPUs with Python 3.8 and PyTorch 1.8. More detailed requirements can also refer to the environment descriptions in our aforementioned source codes.

838 839 840

841 842

843

844 845

846

B DISCUSSION ABOUT RELATED WORK

In this section, we discuss the related literature on machine unlearning, and provide more detailed comparisons of some work with their approaches and motivations.

B.1 MACHINE UNLEARNING

847 Machine unlearning targets to adjust a trained model to scrub the data influence [33, 52, 59]. It is initially proposed to protect data privacy [5, 4, 16], and a series of studies explore probabilistic 848 methods through the differential privacy [16, 21, 47, 57, 50]. Although having the provable guarantee 849 on the unlearning errors, the strong algorithmic assumptions hinders the practical effectiveness [30]. 850 Current research [17, 55, 54, 11, 6, 60, 14, 64] focus more on developing more effective and efficient 851 unlearning methods to approximate the Retrained model, with the given trained model. As for the 852 assumption on data generation, prior works [17, 58, 30, 6] mainly consider all matched forgetting 853 targets, with similar features on the original training tasks. As for the assumption on label generation, 854 most prior works [4, 20, 54, 30, 11, 12] assume the accessibility on the fully identified forgetting 855 dataset, and the complementary is the remaining dataset. One recent work [61] considers unlearning 856 with only a few forgetting samples but requires another generative model to generate approximated data. Our work considers a more practical scenario in which we can conduct mismatched forgetting 858 and use limited identified forgetting data with the remaining set. More details and the intuition of 859 related baseline methods are introduced in Appendix C.

860 861

862

B.2 POSITIVE-UNLABELED LEARNING

Positive-unlabeled learning [10, 45] tries to learn a binary classifier from a few labeled positive samples with the rest unlabeled ones. A series of PU algorithms [40, 8, 9] are developed to train

an accurate binary classifier, and can be roughly divided into two categories [1]. The first branch is cost-sensitive learning, which is related to importance weighting [41]. Given the estimated class prior, these methods [9, 32, 7] can develop an unbiased or consistent risk estimator for PU learning. Another branch of PU learning adopts two heuristic steps to perform binary classification. Such methods [40, 62] first identify reliable negative and positive examples from the unlabeled data, and then conduct semi-supervised learning. The model trained using cost-sensitive learning can also be a recognizer for positive or negative samples [27]. Different from PU learning focusing on binary classification tasks, our work tries to enable more practical scenarios in class-wise unlearning [52] where the class labels and target concepts are decoupled, and we consider the label domain mismatch.

C DETAILS ABOUT CONSIDERED BASELINES AND METRICS

In this section, we provide details about the considered representative baselines for machine unlearning methods, as well as their general intuitions with specific objectives. For the specific hyperparameters adopted in different methods, we keep the same setting with previous related works [30, 11], and the specific values are listed in detail in our source codes. In addition, we introduce the evaluation metrics in detail, corresponding to the implementations in different unlearning scenarios.

C.1 UNLEARNING METHODS

Finetune (FT). Utilizing the catastrophic forgetting [31] in the model (e.g., existed in the continual learning), FT [58] fine-tunes the given trained model partially on \mathcal{D}_r with few training epochs to obtain the θ_{un}^* with the following objective function,

$$L_{\rm FT} = \frac{1}{|\mathcal{D}_{\rm r}|} \sum_{(x,y)\sim\mathcal{D}_{\rm r}} \ell(f(x),y).$$
(8)

Gradient Ascent (GA). Different from the normal gradient descent, GA reverses the gradient signal on \mathcal{D}_{f} to conduct maximization with ascended gradients, resulting in the increasing loss of the forgetting data to obtain the θ_{un}^{*} . The objective is given as follows,

$$L_{\rm GA} = -\frac{1}{|\mathcal{D}_{\rm f}|} \sum_{(x,y)\sim\mathcal{D}_{\rm f}} \ell(f(x),y).$$
(9)

With reverse optimization to maximize the loss on the specific data, the model can approximate θ^* by directly forgetting the learned knowledge represented by the forgetting data.

Random Label (RL). Similar to GA, RL [17] assign the random labels Y^* on the forgetting data in \mathcal{D}_f and fine-tune the given model with it to obtain the unlearned model θ_{un}^* ,

$$L_{\mathrm{RL}} = \frac{1}{|\mathcal{D}_{\mathrm{f}}|} \sum_{(x,y)\sim\mathcal{D}_{\mathrm{f}}} \ell(f(x), y^*).$$

$$\tag{10}$$

Instead of using the original training label on the forgetting data in \mathcal{D}_f , RL can destroy the learned feature by using the random label y^* on \mathcal{D}_f , which violate the minimized loss value.

Influence Unlearning (IU). IU adopts the influence function [33] to estimate the change if the training point is removed from the training loss. It is designed for random data unlearning [52] with the provable guarantee on the unlearning effects. In general, IU estimates the change in model parameters of $\theta_{un}^* - \theta$ and adds the weight perturbation to the given model to obtain the unlearned one. However, it usually requires additional model information and training assumptions for the theoretical guarantee and may suffer hyperparameter tuning with inaccurate hessian estimation [30, 11].

Boundary Shrink (BS). BS [6] is recently proposed for class-wise unlearning, especially on the all matched forgetting. It focuses on the decision spaces [18] of the given trained model. The critical idea is to shift the original decision boundary to imitate the decision behavior of the model retrained from scratch. Motivated by adversarial attacks [43], it proposes a neighbor searching method to identify the nearest but incorrect class labels y_{near} for D_f to guide the model to unlearn the existing class and shift the decision boundary. Using the adversarial attack to find the nearest incorrect label,
 the objective of BS can be formulated as follows,

$$L_{\rm BS} = \frac{1}{|\mathcal{D}_{\rm f}|} \sum_{(x,y)\sim\mathcal{D}_{\rm f}} \ell(f(x), y_{\rm near}),\tag{11}$$

where y_{near} is obtained by first perturbing the forgetting data and getting the newly predicted result as,

x'

$$= x + \epsilon \cdot \operatorname{sign}(\nabla \ell(f(x), y))$$

$$y_{\text{near}} \leftarrow \operatorname{softmax}(f(x'))$$
(12)

926 927

921 922 923

924

925

928

929

930

931

936 937

938

939

940

941

942 943

949 950

966

 L_1 -sparse. Developed based on the conventional FT, L_1 -sparse [30] investigate the model sparsity on machine unlearning. It figures out that model sparsification can benefit the unlearning performance on different perspectives via first pruning and then conducting unlearning. By carrying out pruning and unlearning simultaneously, L_1 -sparse proposes the sparsity-aware unlearning utilizing the L_1 norm-based penalty. The objective is as follows with a hyperparameter γ ,

$$L_{L_1-\text{sparse}} = \frac{1}{|\mathcal{D}_{\mathbf{r}}|} \sum_{(x,y)\sim\mathcal{D}_{\mathbf{r}}} \ell(f(x),y) + \gamma ||\theta^*||, \tag{13}$$

and the general sparsity-aware penalty can also be added to different unlearning methods. In this work, we mainly compare the L_1 -sparse FT as the previous work [30, 11] considered.

SalUn. With the concern on unlearning stability and cross-domain applicability, SalUn [11] introduces the concept of weight saliency in machine unlearning. This innovation directs the attention of unlearning into specific model weights for specific data that need to be unlearned. In general, it first generates the gradient-based weight saliency map inspired by model sparsification [30] with gradient-value thresholding, where the specific generation method is defined as,

$$\mathbf{m}_s = \mathbf{1}(|\nabla_{\theta}\ell(\theta;\mathcal{D}_f)|_{\theta} = \theta_o| \ge \gamma), \quad \theta_u = \mathbf{m}_s \odot (\delta\theta + \theta_o) + (1 - \mathbf{m}) \odot \theta_o, \tag{14}$$

in which $\mathbf{1}(g \ge \gamma)$ is an element-wise indicator function that yields a value of 1 for the i-th element if and 0 otherwise, $|\cdot|$ is an element-wise absolute value operation, and $\gamma > 0$ is a hard threshold. and then conducts saliency-based unlearning using the generated saliency map. Specifically, SalUn adopts RL [17] to fine-tune the forgetting data in $\mathcal{D}_{\rm f}$ on the salience map, and the extended objective is given as follows,

$$L_{\text{SalUn}} = \frac{1}{|\mathcal{D}_{\mathbf{f}}|} \sum_{(x,y)\sim\mathcal{D}_{\mathbf{f}}} \ell_{\theta_u}(f(x), y^*) + \alpha \frac{1}{|\mathcal{D}_{\mathbf{r}}|} \sum_{(x,y)\sim\mathcal{D}_{\mathbf{r}}} \ell(f(x), y), \tag{15}$$

More detailed operations can refer to [11], and we keep the same hyperparameter used in [11] to conduct the class-wise unlearning tasks.

SCRUB. SCRUB is a newly proposed unlearning algorithm based on a novel casting of the problem into a teacher-student framework [37]. It is designed to meet the desiderata of unlearning: efficiently forgetting without hurting the model utility. As the general target of SCRUB in forgetting is application-dependent, it is proposed with a recipe that works across applications: SCRUB is first to strive for maximal forget error, which is desirable in some scenarios like removing bias or restricted contents but not in others like user privacy protection. To address the latter case, SCRUB is integrated with a rewinding procedure that can reduce the forget set error appropriately when required.

Given the original model θ^o as the teacher model, the goal of SCRUB is formatting as training a student model θ^u that selectively obeys the teacher. The overall objective can be divided into two folds, the first is to remember \mathcal{D}_r under the teacher model's guide while the second is to forget \mathcal{D}_f by disobeying the teacher model's guide. To measure the degree to which the student model obeys the teacher model, SCRUB utilizes the following distance measure,

$$d(x;\theta^u) = D_{\mathrm{KL}}(p(f(x;\theta^o))||p(f(x;\theta^u))), \tag{16}$$

where $D_{\rm KL}$ is the KL-divergence and the overall measures of the distance between the student model's and teacher model's prediction distribution. With the aforementioned distance, the objective of SCRUB is as follows,

970
971
$$L_{\text{SCRUB}} = \min_{\theta^u} \frac{\alpha}{N_r} \sum_{x_r \in \mathcal{D}_r} d(x_r; \theta^u) + \frac{\gamma}{N_r} \sum_{(x_r, y_r) \in \mathcal{D}_r} \ell(f(x_r; \theta^u), y_r) - \frac{1}{N_f} \sum_{x_f \in \mathcal{D}_f} d(x_f; \theta^u), \quad (17)$$

972 where the first two parts can be regarded as a variant of distillation from a teacher model on \mathcal{D}_r and 973 the third part is encouraging the student model to disobey the teacher model to forget the target data. 974

Due to the objective design and implementation tricks adopted in SCRUB, we find it may fail to 975 conduct effectively unlearning when the forgetting target consists of a large amount of data (e.g., 976 refer to its good results on Tiny-ImageNet or CIFAR-100 cases compared with that on CIFAR-977 10). We further debugged the failures in Tabel 4. We conjecture it may be due to its specific 978 objective of requiring different prediction results from the original pre-trained model. As there is no 979 hyperparameter that directly controls the forgetting part, the method needs further adjustment when 980 being adopted in unlearning a relatively larger amount (e.g., 4500 samples in one class of CIFAR-10) 981 of target than the original experiments (e.g., 25 or 100 samples) conducted in the paper [37], e.g., 982 sometimes need using lower learning rate for avoiding excessive forgetting.

Table 4: Unlearning using SCRUB with different hyperparameters in the all matched forgetting task on CIFAR-10. $*L = \alpha \cdot d_{kl}(x_r; \theta^u) + \gamma \cdot \ell(f(x_r), y_r) - d_{kl}(x_f; \theta^u)$, recommending with γ =0.99, α =0.001 [37] and lr = 0.01.

Method	Hyperparameter	UA	RA	TA	MIA	Gap↓
Retrained TARF (ours)		0.00 0.00	99.51 98.23	94.69 91.95	$100.00 \\ 100.00$	- 1.01
	γ =0.99 α =0.001	0.00	12.92	12.92	0.00	67.09
	$\gamma = 0.99 \alpha = 0.01$	0.00	18.30	18.52	0.00	64.35
SCDUD*	$\gamma = 0.99 \alpha = 0.1$	0.00	16.99	17.07	0.00	65.06
SCRUD*	omit $-d_{kl}(x_f; \theta^u)$	41.51	99.87	94.44	99.91	10.55
	$\gamma = 0.99 \alpha = 1.0$	0.00	13.71	13.23	99.95	41.82
	$\gamma = 0.99 \ \alpha = 10.0$	0.00	20.18	20.04	100.00	38.46

C.2 EVALUATION METRICS REGARDING DIFFERENT SCENARIOS

In this part, we summarize the following list and tables of the evaluation metrics (adopted from the previous work [30, 11]) and the used labels in different unlearning scenarios,

- Unlearning Accuracy (UA): the accuracy of the unlearned model θ_u on the dataset of target concept $D_{\rm t}$.
- Retaining Accuracy (**RA**): the accuracy of the unlearned model θ_u on retaining dataset D_r .
- Testing Accuracy (**TA**): the accuracy of the unlearned model $\theta_{\rm u}$ on test dataset $D_{\rm test}$ excluding the data belonging to the target concept to be forgotten.
- Model Inversion Attack (MIA): the MIA success rate by a confidence-based MIA predictor of the model θ_u on the dataset of target concept D_t . We follow [30] to implement it to find how many samples in $D_{\rm t}$ can be correctly predicted as a non-training sample by the MIA predictor against $\theta_{\rm u}$. First, we sample a balanced dataset from the retaining dataset $D_{\rm r}$ and the test dataset excluding the forgetting data to train the MIA predictor, then it is used to count the rate of true negative predictions for forgetting data of the target concept.
- 1011 1012 1013

983

984

985

986 987

997 998

999

1000

1002

1004

1008

1009

1010

Generally, in the evaluation phase, we adopt the same labels used in pre-training to measure the 1014 unlearned model. Note that in the model mismatch forgetting, as the model is trained with superclass 1015 labels, the UA is also calculated using the superclass label. Hence, the UA of the Retrained reference 1016 is not equal to 0 as indicated in Table 2, and we compare the methods mainly on the averaged 1017 performance "Gap" (calculated based on the previous four metrics) to the Retrained reference.

1018 In Table 5, we summarize the specific label used in different unlearning scenarios. To provide an 1019 intuitive example that corresponds to the instantiated unlearning tasks like Figure 1, we present 1020 Table 6 to give overall information about the data and labels considered in each metric. 1021

1022

FULL DISCUSSION ABOUT LABEL DOMAIN MISMATCH D 1023

1024

In this section, we discuss the full scenarios of label domain mismatch in class-wise unlearning [58, 1025 17, 6, 30, 11]. Specifically, we will start by why focusing on class-wise unlearning, and then discuss

Used Label	All matched	1 Target mism	atch	Model mismatch		Data mismatch	
UA	Class Labe	l Class Lab	el	Superclass Label	Sup	perclass Label	
RA	Class Label	l Class Lab	el	Superclass Label	Sup	erclass Label	
ТА	Class Label	l Class Lab	el	Superclass Label	Sup	erclass Labe	
MIA	Class Labe	l Class Lab	el	Superclass Label	Sup	erclass Labe	
Table 6: The ev	valuation data ((label number) o	f diffe	erent forgetting scena	arios	with CIFAR-	
Table 6: The ev Data (classes number)	valuation data ((label number) o	f diffe	Model mismatch	arios	with CIFAR-	
Table 6: The ev Data (classes number) UA (D _t)	All matched	Iabel number) o Target mismatch "boy", "girl", "man", "woman", "baby" (5)	f diffe	Model mismatch rt of "people" (1), which is da " and "girl" but with superclas	arios v nta ss label	with CIFAR-	
Table 6: The ev Data (classes number) UA (D _t) RA (D _r)	All matched "boy", "girl" (2) Other classes (98)	Iabel number) o Target mismatch "boy", "girl", "man", "woman", "baby" (5) Other classes (95)	f diffe	Model mismatch rt of "people" (1), which is da " and "girl" but with superclas other part of "people" (1) with the rest superclasses (19)	arios v ata ss label	with CIFAR- Data mismat "people" (1 Other superclasse	
Table 6: The ev Data (classes number) UA (D _t) RA (D _r) TA (D _{test})	All matched "boy", "girl" (2) Other classes (98) Other classes (98)	Iabel number) o Target mismatch "boy", "girl", "man", "woman", "baby" (5) Other classes (95) Other classes (95)	f diffe	Model mismatch rt of "people" (1), which is da " and "girl" but with superclar other part of "people" (1) with the rest superclasses (19) other part of "people" (1) with the rest superclasses (19)	arios v ata ss label 1	with CIFAR- Data mismat "people" (1 Other superclasse Other superclasse	

Table 5: The label used in evaluation metrics on different forgetting scenarios



The left panel shows an example of two-layer label domains; The middle panel is the Venn diagram to show the hierarchical relation; The right panel illustrates the potentials of three critical class-wise unlearning aspects.

Figure 8: Label domain mismatch with the two-layer illustration.

Table 7: Mismatching in the label domain of three critical aspects with a two-layer label structure.

No.	Forgetting data	Model output	Target concept	Comment
1	Class label	Class label	Class label	All matched
2	Class label	Class label	Superclass	Target mismatch
3	Class label	Superclass	Class label	Model mismatch
4	Class label	Superclass	Superclass	Data mismatch
5	Superclass	Class label	Class label	Impractical since $\mathcal{L}_D \succ \mathcal{L}_T$
6	Superclass	Class label	Superclass	Similar to all matched
7	Superclass	Superclass	Class label	Impractical since $\mathcal{L}_D \succ \mathcal{L}_T$
8	Superclass	Superclass	Superclass	All matched

the motivation for investigating its label domain mismatch, with the newly introduced setting being friendly for empirical analysis and further research. In addition, we provide detailed information on our instantiated four tasks using the benchmarked datasets [35]. Finally, we discuss the commonalities of mismatch forgetting scenarios and the general principle of unified framework design.

To begin with, machine unlearning [5, 55, 59, 52] is originally proposed in response to "the right to be forgotten" to protect the data privacy, and recently deep machine unlearning is a timely research topic associated with foundation models which use massive of data to train [37, 3]. The ensuing data regulation concerns have also expanded the original privacy-protecting goal to more general needs and scenarios [60, 44, 14]. As stated in [51, 37, 30], unlearning a subset of the training set has received increasing attention (like removing sensitive information, and inappropriate content). However, the previous scenarios mainly consider the coinciding class labels with the target concept to be unlearned. Although achieving promising results in forgetting, it is still not enough in practice.

1088 Considering the problem setups of unlearning, we have three critical aspects, e.g., the well-trained 1089 machine learning model θ , and the reported data \mathcal{D}_{f} to be unlearned, as well as the target concept. 1090 In previous studies, the three aspects are mainly considered to be under the same label taxonomy. In other words, the unlearning tasks are aligned with the pre-training task, where the latter trains 1091 a multi-class classification model, and the former aims to unlearn a training class. However, in practice, the unlearning request may violate the taxonomy of the pre-training tasks, while the 1093 specific target concepts always exhibit a unified property for specific forgetting data. It naturally 1094 motivates us to consider different label domains of the three aspects of unlearning. As listed 1095 in Table 8, the label domain of data \mathcal{L}_D , the label domain of model output \mathcal{L}_M , and the label domain of target concept \mathcal{L}_T . To begin with, we introduce the relations between two label domains, i.e., \mathcal{L}_1 matches \mathcal{L}_2 ($\mathcal{L}_1 = \mathcal{L}_2$), \mathcal{L}_1 is the subclass domain of \mathcal{L}_2 ($\mathcal{L}_1 \prec \mathcal{L}_2$)³ and \mathcal{L}_1 is the superclass domain of \mathcal{L}_2 ($\mathcal{L}_1 \succ \mathcal{L}_2$)⁴, and we have a practical assumption on the relation between 1099 label domains of forgetting data and target concept, i.e., $\mathcal{L}_D \preceq \mathcal{L}_T$, indicating that the reported 1100 forgetting data should be included in the target concept (as intuitively illustrated in the middle 1101 panel of Figure 8). Considering $\mathcal{L}_D = \mathcal{L}_T$, we can have two possibilities on \mathcal{L}_M , e.g., $\mathcal{L}_M = \mathcal{L}_T$ and $\mathcal{L}_M \neq \mathcal{L}_T$, where the former is regarded as all matched when $\mathcal{L}_D = \mathcal{L}_M = \mathcal{L}_T$ and the 1102 latter is the model mismatch. To be more specific, we consider model mismatch forgetting as 1103 $\mathcal{L}_D = \mathcal{L}_T \prec \mathcal{L}_M$, since $\mathcal{L}_M \prec \mathcal{L}_T$ will have no additional effects on the unlearning when $\mathcal{L}_D = \mathcal{L}_T$ 1104 and we can regard it as similar to the all matched case. Considering $\mathcal{L}_D \prec \mathcal{L}_T$, we can have 1105 the target mismatch forgetting when $\mathcal{L}_D = \mathcal{L}_M$ and data mismatch forgetting when $\mathcal{L}_M = \mathcal{L}_T$. 1106

1107

We summarize the mainly considered mismatch
cases in Table 8, which can serve as a general
reference for further research on constructing the
unlearning tasks. In the following, we further
explain the procedure of task instantiating and discuss the other potential scenarios with the typical
two-layer label structure considered in the main
text and an additional three-layer label structure.

Label Domain $\mathcal L$	Relation of Da	ta \mathcal{L}_D	, Mode	\mathcal{L}_M	, and T	larget \mathcal{L}_T
All matched	$ $ \mathcal{L}_D	=	\mathcal{L}_T	=	\mathcal{L}_M	
Target mismatch	\mathcal{L}_M	=	\mathcal{L}_D	\prec	\mathcal{L}_T	
Model mismatch	\mathcal{L}_D	=	\mathcal{L}_T	\prec	\mathcal{L}_M	
Data mismatch	$ $ \mathcal{L}_D	\prec	\mathcal{L}_T	=	\mathcal{L}_M	

Table 8: considering **label domain** relations of three critical aspects in class-wise unlearning.

1115 1116

1117

D.1 A TWO-LAYER LABEL STRUCTURE OF MISMATCH

In Figure 8, we first show the illustration of a two-layer label structure and the three aspects of unlearning, i.e., forgetting data, model outputs, and target concept. Without losing generality, we utilize the class labels and superclass information (refer to the official information in CIFAR-100 [35]) for consideration. Then we have a two-layer label structure representing different knowledge regions.

1122 Given two potential label domains in each aspect, we can totally get the 8 scenarios list in Table 7. 1123 The first 4 scenarios are mainly considered and detailedly introduced in the main text. For the 1124 rest 4 scenarios (i.e., No. 5-8), we consider some (i.e., No. 5 and No. 7) to be impractical as the 1125 label domain of forgetting data is larger than the target concept, which means that the unlearning requests identify more forgetting data than the true target concept. It should be more reasonable that 1126 only limited forgetting data are identified by server users or internal examiner [34] in real-world 1127 applications. Therefore, we mainly consider the forgetting data \mathcal{D}_{f} belongs a part of or equals to 1128 the overall data \mathcal{D}_t of the target concept. As for No. 6 and No. 8 cases, the former is similar to the 1129 conventional all matched forgetting since the forgetting data has the same label domains with the 1130

¹¹³¹ ${}^{3}\mathcal{L}_{1} \prec \mathcal{L}_{2}$: For any label $y \in \mathcal{L}_{1}$, there exist label $y' \in \mathcal{L}_{2}$ that instance being labeled with y can also being labeled with y', but not all instance being labeled with y' can be labeled with y.

¹¹³³ ${}^{4}\mathcal{L}_{1} \succ \mathcal{L}_{2}$: For any label $y \in \mathcal{L}_{2}$, there exist label $y' \in \mathcal{L}_{1}$ that instance being labeled with y can also being labeled with y', but not all instance being labeled with y' can be labeled with y.



To be more specific, there are two groups of cases in the third part. For No. 5, 6, and 7, since they also can be represented using a two-layer structure, the forgetting dynamics are similar to that in target, model, and data mismatch forgetting. By contrast, in No. 8, 9, and 16, all three label domains exist in the three aspects of class-wise unlearning, which is worthy of further discussion.

- 1173
- 1174 D.3 FURTHER EXPLORATION ON THE OTHER 6 DIFFERENT SCENARIOS

In this part, we further discuss the 6 different scenarios discovered by constructing the three-layer
 label structure. We illustrated these forgetting tasks in Figure 10 and discuss them as follows,

1178 - No. 5&16 In the two scenarios, the model output has the most fine-grained label domain (e.g., 1179 sub-set as illustrated in Figure 9) for representation. At the same time, the target concept is broader than both model output and identified forgetting data. Different from the aforementioned target 1180 mismatch, the mismatch degree of this task is larger (e.g., superclass level) than the previous one (e.g., 1181 class level). In other words, the model output further loses the entanglement of feature representation 1182 of the samples belonging to target concept (compared with the original setups of target mismatch). 1183 To simulate the case, we employ the same model pre-trained by class in target mismatch, but enlarge 1184 the target concept (consists of 7 classes with similar semantic features, instead of the original 5) and 1185 change the forgetting data (2 class as the given forgetting data in No.5 and 3 classes in No.16). 1186

1187 - No. 8&7. Similar to the previous No. 5, the target concept in these tasks is also broader than the label domains of the identified forgetting data. However, in these two scenarios, the model output is

1190					~
1191	No.	Forgetting data	Model output	Target concept	Comment
1192	1	Sub-set	Sub-set	Sub-set	All matched
1193	2	Sub-set	Sub-set	Class label	Target mismatch
1194	3	Sub-set	Class label	Sub-set	Model mismatch
1105	4	Sub-set	Sub set	Class label	Data mismatch
1100	6	Sub-set	Superclass	Sub-set	Different
1190	7	Sub-set	Superclass	Superclass	Different
1197	8	Sub-set	Class label	Superclass	Different
1198	9	Sub-set	Superclass	Class label	Different
1199	10	Class label	Class label	Class label	All matched
1200	11	Class label	Class label	Superclass	Target mismatch
1201	12	Class label	Superclass	Class label	Model mismatch
1202	13	Class label	Superclass	Superclass	Data mismatch
1203	14	Class label	Sub-set	Sub-set	Impractical since $\mathcal{L}_D \succ \mathcal{L}_T$
1204	15	Class label	Sub-set	Class label	Similar to all matched
1205	10 17	Class label	Sub-set	Superclass	Impractical since $C_{-} \subseteq C_{-}$
1200	18	Class label	Superclass	Sub-set	Impractical since $\mathcal{L}_D \succ \mathcal{L}_T$
1200	10	Cumarala a	Cupar-1	Sub Set	
1207	19 20	Superclass	Superclass Class label	Superclass Class label	All matched Impractical since $C = \subseteq C$
1208	20	Superclass	Class label	Superclass	Similar to all matched
1209	22	Superclass	Superclass	Class label	Impractical since $\mathcal{L}_D \succ \mathcal{L}_T$
1210	23	Superclass	Sub-set	Sub-set	Impractical since $\mathcal{L}_D \succ \mathcal{L}_T$
1211	24	Superclass	Sub-set	Class label	Impractical since $\mathcal{L}_D \succ \mathcal{L}_T$
1212	25	Superclass	Sub-set	Superclass	Similar to all matched
1213	26	Superclass	Class label	Sub-set	Impractical since $\mathcal{L}_D \succ \mathcal{L}_T$
1214	27	Superclass	Superclass	Sub-set	Impractical since $\mathcal{L}_D \succ \mathcal{L}_T$
1015					
1016		Class Representation	on Cla	ss Representation	Class Representation
1216		Class Representation	on Cla	ss Representation	Class Representation
1216 1217		Class Representatio	on Cla	ss Representation	Class Representation
1216 1217 1218		Class Representation	on Cla	ss Representation	Class Representation
1216 1217 1218 1219		Class Representation		ss Representation	Class Representation
1216 1217 1218 1219 1220		Class Representation		ss Representation	Class Representation
1216 1217 1218 1219 1220 1221		Class Representation	on Clar	ss Representation	Class Representation
1216 1217 1218 1219 1220 1221 1222		Class Representation		ss Representation	Class Representation
1216 1217 1218 1219 1220 1221 1222 1223		Class Representation		ss Representation	Class Representation
1216 1217 1218 1219 1220 1221 1222 1223 1224		Class Representation		ss Representation	Class Representation
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225		Class Representation		ss Representation	Class Representation
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226		Class Representation	n Clas	ss Representation	Class Representation View Provident View Pr
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227		Class Representation	n Clas	ss Representation	Class Representation Figure 1 Figure
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227		Class Representation	n Clas	ss Representation	Class Representation View Presentation View Presentation No. 6 Class Representation
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228		Class Representation	n Clas	ss Representation	Class Representation View Presentation No. 6 Class Representation
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228 1229		Class Representation	n Clas	ss Representation	Class Representation Final Action Acti
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230		Class Representation	n Clas	s Representation	Class Representation Final Action Acti
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231		Class Representation	n Clas	s Representation	Class Representation Figure 1 Figure
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232		Class Representation	n Clas	ss Representation	Class Representation No. 6 Class Representation Class Representation
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232		Class Representation	n Clas	ss Representation	Class Representation No. 6 Class Representation Class Representation
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234		Class Representation	n Clas	s Representation	Class RepresentationImage: colspan="2">Image: colspan="2" Colspan="2">Image: colspan="2" Colspa
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235		Class Representation	n Clas	s Representation	Class Representation Final Property State Property Pr
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236		Class Representation	n Clas	s Representation	Class RepresentationImage: colspan="2">Image: colspan="2">Class RepresentationImage: colspan="2">Image: colspan="2">Image: colspan="2">Class RepresentationImage: colspan="2">Image: colspan="2">Image: colspan="2">Class RepresentationImage: colspan="2">Image: colspan="2">Image: colspan="2"Image: colspan="2">Image: colspan="2" </td

1188Table 9: Mismatching in the label domain of three critical aspects with a three-layer label structure.1189

Figure 10: Illustration of 6 scenarios different from the four major tasks according to Table 9.

varied which controls the entanglement of target samples. To construct these two forgetting tasks, we
respectively adopt the models pre-trained by class labels and superclass, and use the same forgetting data (with 1 class) to investigate the performance change using our TARF and other baselines.

18.44

92.21

96.06

98.43

78.20

82.32

CIFAR-100	Metric	UA	RA	TA	MIA	Gap↓	Metric	UA	RA	TA	MIA	Gap
Retrained		0.00	97.85	73.72	100.00	-		0.00	97.85	73.72	100.00	-
FT [58]	1	67.52	96.43	72.96	41.14	32.72	1	53.11	94.64	71.23	52.70	26.
RL [56]	No. 5	68.57	96.12	72.58	41.17	33.15	No. 16	53.90	96.94	73.07	53.56	25.
GA [28]		38.03	97.00	70.98	76.92	16.75		32.24	95.73	69.99	77.62	15
TARF (ours)		0.00	96.58	72.03	100.00	0.74		0.00	96.98	72.87	100.00	0.4
Retrained		0.00	97.85	73.72	100.00	-		0.00	98.50	80.15	100.00	-
FT [58]		74.09	97.19	74.01	36.71	34.58		95.16	94.98	78.68	13.06	46
RL [56]	No. 8	76.04	96.76	72.88	36.00	35.49	No. 7	91.51	96.98	80.11	47.24	36
GA [28]		49.47	98.92	72.94	77.96	18.34		15.91	98.64	80.27	93.82	5.
TARF (ours)		0.00	96.22	72.43	100.00	0.73		0.00	96.54	79.23	100.00	0.
Retrained		88.22	98.52	84.42	22.22	-		88.22	98.58	78.50	25.78	
FT [58]		94.33	95.00	78.77	13.67	5.96		91.78	95.02	78.90	18.44	3.
RL [56]	No. 6	84.22	96.96	80.18	65.77	13.34	No. 9	96.97	70.22	80.24	94.67	26
		10.44	00.00	70.00	02 (7	27.22		10.11	05 07	77 56	01 50	24

92.67

19.17

37.23

2.31

19.11

89.12

95.27

97.23

77.56

79.21

91.56

24.32

34.79

1.11

1242 Table 10: Results (%) of unlearning with different model structures. All methods are trained on the 1243 same backbone, i.e., the basis of unlearning initialization is the same (except for retraining from 1244 scratch). Values are percentages. Bold numbers are superior results. \downarrow indicates smaller are better.

1261 - No. 6&9. In the last two scenarios, the forgetting tasks are more similar to the previous model 1262 mismatch forgetting. However, the distinguishable difference is that the label domain of model 1263 outputs can be much different from the identified forgetting data. In the No. 6 task, the target concept 1264 is aligned with the identified forgetting data, while since the remaining data is more than the original 1265 model mismatch forgetting, the task separation could be harder than the previous. In the No. 9 task, 1266 we can find that it is a complex scenario where the target concept is broader than the forgetting data but included in the same superclass. In both tasks, we use 1 class data as the forgetting data. 1267

1268 To further understand the properties of unlearning in these tasks, we conducted additional experiments 1269 and summarized the results in Table 10. We can find the empirical results well demonstrate the 1270 conceptual conjectures in the previous discussion, and the representative baselines exhibit varied 1271 performance gap with the Retrained reference. Among them, our TARF can consistently achieve the 1272 better performance regarding to the Gap.

1273

1258

1259 1260 GA [28]

TARF (ours)

1274 D.4 SPECIFIC INFORMATION OF THE INSTANTIATED TASKS 1275

1276 For the four major scenarios (i.e., conventional all matched forgetting, target mismatch forgetting, 1277 model mismatch forgetting, and data mismatch forgetting) considered in our work, we provide the 1278 dataset construction and partition details in this section. Note that we focus on class-wise unlearning 1279 in this work, which is different from random data forgetting that uniformly samples the forgetting target of all classes in the training dataset. 1280

1281 Forgetting target. In previous works [58, 6], the target concept to be forgotten is mainly considered 1282 as all matched where $\mathcal{D}_t = \mathcal{D}\{y = y_f\}$ has the same label domains (exactly same labels) with the 1283 pre-training task and forgetting data $\mathcal{D}_f = \mathcal{D}\{y = y_f\}$. In contrast, we assume that the target concept 1284 can be decoupled from the class label in practical unlearning requests. As illustrated in Figure 1, 1285 we further instantiate with three forgetting tasks given $\mathcal{D}_{f} = \mathcal{D}\{y = y_{f}\}$ with the superclass labels \mathcal{Y}' of \mathcal{Y} (classes): i) model mismatch forgetting, e.g., $\mathcal{D}_t = \mathcal{D}\{y = y_t\}$ and $y_t \subseteq y'_f$ where $y'_f \in \mathcal{Y}'$ 1286 given the model trained on \mathcal{Y}' ; ii) target mismatch forgetting, e.g., $\mathcal{D}_t = \mathcal{D}\{y = y'_f\}$ given the model 1287 trained on \mathcal{Y} ; iii) data mismatch forgetting, e.g., $\mathcal{D}_t = \mathcal{D}\{y = y'_t\}$ given the model trained on \mathcal{Y}' . 1288

To ease the research investigation and empirical verification, we adopt the commonly used [37, 30, 1290 11, 12] benchmark CIFAR-10 and CIFAR-100 for constructing the pre-training task for unlearning. 1291 Specifically, the official class labels are kept as classes for ordinary setup, and we provide the superclass information referring to the pre-defined lists [35] of CIFAR-100. Since there is no official superclass information for CIFAR-10 dataset, we manually grouped the classes of CIFAR-10 1293 according to their semantic feature similarity and finalized 5 superclass clusters consisting of 2 classes 1294 in each. The full structured label layers information is summarized in Tables 12 and 13. For all the 1295 unlearning scenarios where the label domain of model output is the superclass, we will first use the



Taking the *CIFAR-100* [35] dataset, we instantiate four unlearning tasks given the same forgetting data with the class labels of "boy" and "girl":
 a) all matched forgetting (conventional scenario): unlearn "boy" and "girl" with the model trained on the classes; b) target mismatch forgetting:
 unlearn "people" with the model trained on the classes; c) model mismatch forgetting: unlearn "boy" and "girl" with the model trained on the superclass; d) data mismatch forgetting: unlearn "people" with the model trained on the superclass.

Figure 11: Illustrations of class representation with the four unlearning scenarios.

superclass information to train the 20-class and 5-class classification models respectively. For the
specific data partition in unlearning requests, we randomly sampled two classes in CIFAR-100 and
one class in CIFAR-10 as forgetting data and kept the setup across the four forgetting tasks as well as
other experiments. For other additional experimental setups, we will state them at the near positions.

Table 11: Basic setup about unlearning scenarios. More illustrations can be found in Appendix D.4.

Dataset Forg	getting Data	Setup	All matched	Model mismatch	Target mismatch	Data mismatch
CIFAR-10 "au	itomobile"	Training Class	10	5	10	5
		Target Concept	"automobile"	"automobile"	"vehicle"	"vehicle"
CIFAR-100 CIFAR-100	ov" "girl"	Training Class	100	20	100	20
		Target Concept	"boy", "girl"	"boy", "girl"	"people"	people"

Table 13: Full list of the 5-class classification on CIFAR-10 with its manually set superclass [35].

Table 14: Specific training set data partition corresponding to four major forgetting tasks.

Superclass (5)	Classes (2 for each superclass)	Forgetting Tasks	Identified	Unidentified
1	airplane, bird	All matched	$\mathcal{D}_{\mathrm{f}} = \mathcal{D}_{\mathrm{t}}$	$\mathcal{D}_{\mathrm{un}} = \mathcal{D}_{\mathrm{r}}$
2	automobile, truck	Target mismatch	$\mathcal{D}_{\mathrm{f}} \subset \mathcal{D}_{\mathrm{t}}$	$\mathcal{D}_{un} = \mathcal{D}_{fr} \cup \mathcal{D}_r$
3	cat, dog	Model mismatch	$\mathcal{D}_{\mathrm{f}} = \mathcal{D}_{\mathrm{t}}$	$\mathcal{D}_{\mathrm{un}}=\mathcal{D}_{\mathrm{r}}$
4	deer, frog	Data mismatch	$\mathcal{D}_{\mathrm{f}} \subset \mathcal{D}_{\mathrm{t}}$	$\mathcal{D}_{un} = \mathcal{D}_{fr} \cup \mathcal{D}_r$
5	horse, ship			-

D.5 DISCUSSION ON THE PRACTICALITY OF LABEL DOMAIN MISMATCH

Machine unlearning is originally proposed in response to data regulations [5, 26], which are primarily motivated by a desire to protect data owners' right to withdraw from the learning process. However, regarding its technical nature of mitigating the data influence from a trained model [59], unlearning is actually given broader significance in the context of trustworthy AI [4,5], like studies for mitigating bias and unfairness [6], addressing safety issues [7], erasing the NSFW generation [8,9]. It is worth noting that these trustworthy requirements may generally exhibit different concerns from the original training tasks. Motivated by the research problem raised in Section 3.1, our work focuses on a critical problem from the assumption view, i.e., the unlearning request may have a different taxonomy from the original tasks, for which we model the three mismatched scenarios for systematical exploration.

1352	\mathbf{C}	$(5, 1, \dots, 1, \dots, 1, \dots)$
1353	Superclass (20)	Classes (5 for each superclass)
1354	aquatic mammals	beaver, dolphin, otter, seal, whale
1355	fish	aquarium fish, flatfish, ray, shark, trout
1356	flowers	orchids, poppies, roses, sunflowers, tulips
1357	food containers	bottles, bowls, cans, cups, plates
1358	fruit and vegetables	apples, mushrooms, oranges, pears, sweet peppers
1359	household electrical devices	clock, computer keyboard, lamp, telephone, television
1360	household furniture	bed, chair, couch, table, wardrobe
1361	insects	bee, beetle, butterfly, caterpillar, cockroach
1362	large carnivores	bear, leopard, lion, tiger, wolf
1363	large man-made outdoor things	bridge, castle, house, road, skyscraper
1364	large natural outdoor scenes	cloud, forest, mountain, plain, sea
1365	arge omnivores and nerbivores	for porcuping possum reason skupk
1266	non insect invertebrates	lobater angil anider worm
1267	non-insect invertebrates	haby boy girl man woman
1000	rentiles	crocodile dinosaur lizard snake turtle
1000	small mammals	hamster mouse rabbit shrew squirrel
1369	trees	manle oak palm nine willow
1370	vehicles 1	bicycle bus motorcycle pickup truck train
1371	vehicles 2	lawn-mower rocket streetcar tank tractor
1372	, enneres 2	aun moner, rocket, streeteur, ank, tractor

Table 12: Full list of the 20-class classification on CIFAR-100 with its official superclass labels [35].

1373

1374

Here, we discuss some practical use cases for the three newly proposed settings. For example, 1) the 1375 label domain mismatch may exist in some recommendation tasks [15] or other generative tasks like 1376 image generation [14] with diversified user feedback (for which we have presented a case study on 1377 concept-forgetting on Stable Diffusion in Appendix E). 2) When the users raise unlearning requests 1378 for some representative disliked item with a message of "don't recommend this kind of thing", it is 1379 similar to target mismatch forgetting. In addition, 3) when debugging the pre-trained model with some 1380 spurious correlation or safety concerns [3, 14, 39], it is similar to model or data mismatch forgetting 1381 as we only have the forgetting cases (e.g., some figures including NSFW content or adversarial features) that may not be aligned with the taxonomy of the pre-training task. We hope our exploration 1382 can provide insights for further consideration of specific practical applications. 1383

1384 1385

1386

D.6 DISCUSSION ON THE SCENARIO COMMONALITIES AND FRAMEWORK PRINCIPLES

The three mismatch scenarios, i.e., target mismatch forgetting, data mismatch forgetting, and model mismatch forgetting, share the common challenge of representation mismatch between the pre-trained model, the identified forgetting data, and the target concept to be forgotten. It breaks the assumption in all-matched scenarios that the three are matched [51, 59, 11, 30] and can result in extra-/ineffectiveforgetting in unlearning tasks, as demonstrated in our Figure 2. Specifically, in the target/data mismatch forgetting, the target concept can be wider than the identified forgetting data; while in the model mismatch forgetting, it can be smaller than the coarse-grained model representation.

1393 To build a unified framework like our TARF, it requires considering the aforementioned two issues, i.e., 1394 insufficient representation (refer to the observation 3.1) in target/data mismatch, and decomposition 1395 lacking (refer to the observation 3.2) in the model mismatch. The former requires a flexible controller 1396 for forgetting strength while and latter requires a simultaneous consideration on forgetting and 1397 retaining. Thus, based on the general equation in Eq. (5), we set two sub-objectives (annealed 1398 forgetting and target-aware retaining) to decompose the learned representation and control the 1399 forgetting strength by the instance-wise weighting mechanism which selects the targeted-aware forgetting data. Then, TARF becomes a unified framework for the three mismatched scenarios. Note 1400 that in our presentation, TARF is illustrated with three phases to better explain its functionality, while 1401 it is not three independent parts but unified in a general objective with specific parameters. 1402

1404 Ε ALGORITHM IMPLEMENTATION AND EXPLANATION 1405

In this section, we present the pseudo-code of our proposed TARF and its variant, as well as additional discussions to enhance the understanding of our methods. Here we summarize the detailed procedure 1408 of algorithm implementation in Algorithm 1 and Algorithm 2. In detail, Algorithm 1 identifies the potential target using the class labels, while Algorithm 2 can use the instance level information.

As introduced in Section 3.3, the objective of our TARF is defined as follows,

$$L_{\text{TARF}} = \underbrace{k(t) \cdot \left(-\frac{1}{|\mathcal{D}_{f}|} \sum_{(x,y) \sim \mathcal{D}_{f}} \ell(f(x), y)\right)}_{\text{Annealed Forgetting } L_{f}(k)} + \underbrace{\frac{1}{|\mathcal{D}_{\text{un}}|} \sum_{(x,y) \sim \mathcal{D}_{\text{un}}} \ell(f(x), y) \cdot \tau(x, y, t)}_{\text{Target-aware Retaining } L_{u}(\tau)}$$
(18)

where k(t) and $\tau(x, y, t)$ are two training-time-related hyperparameters to deal with the mismatch 1416 issues raised in our new settings. Specifically, we set a learning-rate-reduced k(t) as, 1417

> $k(t) = k \cdot (T - t - t_0)/T, \quad t \in [0, T],$ (19)

1419 where T indicates the total training time (e.g., epochs), and the value of k(t) decreases with the 1420 training process. On the other hand, we have the following indicator to measure the model prediction 1421 consistency with the training data $I_{con}(x, y, \theta) = |\ell_{f_{\theta}}(x, y) - \ell_{f_{\theta^*}}(x, y)|$, with which we set $\tau(x, y, t)$ 1422 as follows,

1433

1418

1406

1407

1409

1410

$$\tau(x, y, t) = \begin{cases} 0 & I_{\text{con}}(x, y, \theta_{t_1}) > \beta \text{ or } t < t_1 & \text{*Unconf. Retain,} \\ 1 & I_{\text{con}}(x, y, \theta_{t_1}) < \beta \text{ and } t \ge t_1 & \text{*Conf. Retain,} \end{cases}$$
(20)

where t_1 is a time stamp to control the start of pursuing the retaining part. The overall two dynamic 1426 hyperparameters can divide the whole unlearning process into three phases as illustrated in Figure 4. 1427

1428 In the implementation, the β value is estimated with the ranked loss difference of each class and the 1429 prior information about the target concept to control the forgetting strength in unlearning. For the task feasibility, we will generally assume the amount of false remaining data or classes is known at 1430 our target/data mismatch forgetting, following a similar setup in learning from label noise [22] where 1431 the noise rate can be estimated and utilized as prior information. 1432

Annealed Forgetting. For the forgetting target, we adopt the gradient ascent on the given forgetting 1434 data to unlearn it. However, to approximate the retrained model, the intuition is not to pursue the 1435 maximization of the risk on this part of the data but to destroy the learned feature on the given model. 1436 So we introduce a learning-rate-reduced k(t) to realize the annealed gradient ascent where $t_0 = 1$ 1437 is adopted for target or data mismatch forgetting, and the value of k(t) decreases with the training 1438 process. Resulting in destroyed features, gradient ascent on this part of data also constructs the 1439 dynamic information for differentiating the data of different consistency on its loss values, making 1440 the risks of the false retaining data higher than the rest, and helping to filter retaining data. 1441

1442 Target-aware Retaining. For the retaining part, we need to selectively learn the data from the remaining set, since the complementary dataset may be biased with unidentified forgetting data. 1443 Compared with other remaining data, the false retaining data is easy to be affected by similar feature 1444 representation as indicated in Figure 3(a). Thus, we can have $\tau(x, y, t)$ where we can divide the 1445 remaining set into unconfident/confident parts to note the estimated retaining data like Figure 5(a). 1446 $t_1 = 2$ is adopted at target and data mismatch tasks, and β can be estimated by the prior information 1447 about the specific unlearning request and the rank of loss values. By simultaneously conducting 1448 gradient ascent on forgetting data and selective gradient descent on confident retaining data, we can 1449 better restrict the forgetting region and deconstruct the entangled feature representation (refer to the 1450 middle of Figure 5(b) where we reveal the feature decomposition in deeper layers of model structure 1451 using ResNet). Finally, with the partial objective of retaining, it can approximate the retrained 1452 reference (refer to the right of Figure 5(b)).

1453

1455

To demonstrate the compatibility, we also extend the idea of this work and investigate the performance 1456 of TARF on the specific text-to-image generation task with stable diffusion [14, 11], and presented 1457 the generated images by the original model and unlearned model in Tables 15 and 16.

1458	Algor	ithm 1 TARF
1459	Input	: Training dataset $D = \{(x_i, y_i, s_i)\}_{i=1}^n$, where $s_i = 1$ indicates the identified forgetting dat,
1460	otherv	vise the data is recognized to be unlabeled for unlearning, learning rate η , number of epochs T,
1461	batch	size m, number of batches M, data $x \in \mathcal{X}$, label $y \in \mathcal{Y}$, original trained model θ , loss function
1462	ℓ , init	ialized indicator value τ with the threshold β .
1463	Outp	ut: model θ^T ;
1464	1: f o	r mini-batch = $1, \ldots, M$ do
1465	2:	Sample a mini-batch $\{(x_i, y_i)\}_{i=1}^m$ from D
1466	3:	$\{\ell(x_i, y_i)\}_{i=1}^m \leftarrow \theta. \text{forward}(f_\theta, \{(x_i, y_i)\}_{i=1}^m),$
1467	4:	Collect the initial training accuracy in each class based on $\{\ell(x_i, y_i)\}_{i=1}^m$,
1468	5: e i	nd for
1469	6: f o	$\mathbf{r} = \mathrm{epoch} = 1, \dots, T \mathbf{do}$
1470	7:	Update $k(t)$ according to Eq. 7,
1471	8:	if epoch $< t_0$ then
1472	9:	$ au \leftarrow 0$
1473	10:	
1474	11:	compute β in Eq. / according to the rank of class accuracy difference, and update τ .
1475	12:	end II for mini botch 1 M do
1476	13:	for mini-batch = 1,, M do Somple a mini batch $[(m, u, a)]^m$ from D
1477	14:	Sample a mini-batch $\{(x_i, y_i, s_i)\}_{i=1}$ from D Assign different weights for identified target semples and the rest retaining data
1/170	15:	Assign unreferit weights for identified target samples and the fest retaining data,
1479	16:	$L_{\text{TARF}} = k(t) \cdot \left(-\frac{1}{ \mathcal{D}_{\text{f}} } \sum_{(x,y)\sim\mathcal{D}_{\text{f}}} \ell(f(x),y) \right) + \frac{1}{ \mathcal{D}_{\text{un}} } \sum_{(x,y)\sim\mathcal{D}_{\text{un}}} \ell(f(x),y) \cdot \tau(x,y,t),$
1480	17:	$\theta \leftarrow \theta - \eta \nabla_{\theta} L_{\text{TARF}}(D, D_{\text{f}}, f, \tau)$
1481	18:	end for
1482	19: e i	nd for
1483		

In detail, we aim to unlearn the image generation of a concept with its specific prompt like "a photo of 1485 a tench". To simulate the practical unlearning request (e.g., the user raises the request of unlearning 1486 a specific concept with some identified generation examples, and the developer needs to adjust the 1487 model to forget the concept), we construct the given dataset consisting of limited forgetting data and 1488 the unidentified remaining data for unlearning, which corresponds to the data mismatch forgetting 1489 task. Then we compare the image generation on the original stable diffusion, the unlearned model 1490 with certain label (CL) mismatching [11], and that with our TARF. Note that here we recognize 1491 ESD [14] as a performance upper bound and do not compare it, since it is the same for all matched 1492 settings with fully identified forgetting data (as it directly encourages the model to unlearn the concept 1493 from text semantics). For this exploration of TARF, we adopt the instance-wise identification during the forgetting process as described in Algorithm 2, to unlearn the target concept with the given limited 1494 forgetting data and pursue retaining the selected remaining data with lower loss values. 1495

The results in Tables 15 and 16 demonstrate that our TARF can achieve better forgetting results given the limited identified forgetting data, with proper target identification in the remaining set, while CL using only identified forgetting data can not unlearn the concept well as the generated examples still maintain some semantic features belongs to the target concept (like "tench" or "English springer").

1500

1501 E.2 DISCUSSION ON THE ALGORITHM COMPUTATION COST 1502

We would acknowledge that TARF may require more time in unlearning compared with some methods
like GA, which only uses the forgetting data (which sometimes can be extremely limited than other
retaining data) for unlearning, while those methods may suffer from excessive forgetting and results
in inaccurate unlearning across different scenarios. Regarding the metric "TIME", it originally means
to avoid some methods that consume too much time compared with that of Retrained (Ref.). From
this perspective, these current methods and TARF actually fall in the acceptable time range, and the
efficiency gap between existing explorations in that range is indeed not a bottleneck based on Table 1.

From the methodology perspective, the three separately presented phases are integrated in a unified framework, instead of adding extra phases before and after Phase II. Compared to other approximate unlearning methods, the unique operation is target identification by comparing the output information

1512	Alg	orithm 2 TARF-I: generalized version on instance-wise identification
1513	Inp	ut: Training dataset $D = \{(x_i, y_i, s_i)\}_{i=1}^n$, where $s_i = 1$ indicates the identified forgetting dat,
1514	othe	rwise the data is recognized to be unlabeled for unlearning, learning rate η , number of epochs T,
1515	batc	h size m, number of batches M, data $x \in \mathcal{X}$, label $y \in \mathcal{Y}$, original trained model θ , loss function
1516	ℓ , in	itialized indicator value τ with the threshold β .
1517	Out	put: model θ^T ;
1518	1:	for mini-batch = $1, \ldots, M$ do
1519	2:	Sample a mini-batch $\{(x_i, y_i)\}_{i=1}^m$ from D
1520	3:	$\{\ell(x_i, y_i)\}_{i=1}^m \leftarrow \theta. \text{forward}(f_\theta, \{(x_i, y_i)\}_{i=1}^m),$
1521	4:	Collect the initial loss values in each training samples based on $\{\ell(x_i, y_i)\}_{i=1}^m$,
1522	5:	end for
1523	6:	for epoch $= 1, \ldots, T$ do
1524	7:	Update $k(t)$ according to Eq. 7,
1525	8:	if epoch $< t_0$ then
1526	9:	$ au \leftarrow 0$
1527	10:	else
1500	11:	compute β in Eq. 7 according to the rank of difference in instance loss values, and update τ .
1520	12:	end if
1529	13:	for mini-batch = $1, \ldots, M$ do
1530	14:	Sample a mini-batch $\{(x_i, y_i, s_i)\}_{i=1}^m$ from D
1531	15:	Assign different weights for identified target samples and the rest retaining data,
1532	16.	$I_{\pi,\pi\pi} = k(t) \cdot \left(-\frac{1}{2} \sum_{n=1}^{\infty} \ell(f(x), y) \right) + \frac{1}{2} \sum_{n=1}^{\infty} \ell(f(x), y) \cdot \tau(x, y, t)$
1533	10.	$L_{\text{TARF}} = \kappa(\iota) \cdot \left(-\frac{ \mathcal{D}_{f} }{ \mathcal{D}_{f} } \sum_{(x,y) \sim \mathcal{D}_{f}} \iota(f(x), g) \right) + \frac{ \mathcal{D}_{un} }{ \mathcal{D}_{un} } \sum_{(x,y) \sim \mathcal{D}_{un}} \iota(f(x), g) \cdot f(x, g, \iota),$
1534	17:	$\theta \leftarrow \theta - \eta \nabla_{\theta} L_{\text{TARF}}(D, D_{\text{f}}, f, \tau)$
1535	18:	end for
1536	19:	end for
1537		
1538		

of the unlearned model with the original model for weight assignment, which has similar or less
computation than other advanced designs that consider the sparse regularization [30] or compute the
gradient mask for the original model [11].

E.3 DISCUSSION ON TARF WITH LIMITED CLASS INFORMATION

The phase 1 of TARF is for target identification in the target mismatch forgetting where the target concept is wider than the given forgetting data (e.g., forgetting "people" given "boy" and "girl"). The class information may affect the accurate identification of the target concept but not the rationality of our framework. In our experimental setup, the class information is available in TARF as the class labels are used in pre-training, while the information of the target concept is given by the number of extra classes instead of the superclass label. Regarding the unavailable or implicit class information, first, if the class (i.e., the subclasses w.r.t. target concept) is not available, TARF may also utilize the model prediction to obtain the pseudo labels to conduct the task; Second, if the extra forgetting target beyond the identified data is not restricted as classes, it may require that given forgetting data can well represent the target concept (e.g., the false retaining data should be easier affected than the other retaining data). We acknowledge that both scenarios would lead to a larger performance gap with Retrained reference, as it is a generally more challenging scenario affecting the task achievability to all of the approximate unlearning methods. We believe it worth future effort to explore.

Table 15: Image generation results of unlearned Stable Diffusion in the Data mismatch forgetting, compared with the original stable diffusion, certain label (CL) unlearning [11], and our TARF. The specific prompt used in the image generation is "a photo of tench".

1571			
1572	Original	Unlearned	Unlearned
1573	Stable Diffusion	by CL [11]	by TARF
1574			
1575		- NO REC	3 6 120
1576			
1577			
1578			
1579	(MIII)		
1580			
1581		AX	
1582			
1583			Sec
1584			
1585			
1586			
1587			
1588			
1589		and and	Star
1590		6-16 24	
1591		CEAN AZ	SEAN AND
1592			
1593		GUT	a la
1594			
1595			
1590			
1597			
1590			
1600	THE		
1601			
1602			
1603			
1604			
1605			
1606			
1607			
1608			
1609			
1610			
1611			
1612			Maria Maria
1613			
1614	Contraction of the second seco		
1615			
1616			
1617	2		
1618			

1622Table 16: Image generation results of unlearned Stable Diffusion in the Data mismatch forgetting,1623compared with the original stable diffusion, certain label (CL) unlearning [11], and our TARF. The1624specific prompt used in the image generation is "a photo of English springer".

1625			
1626	Original	Unlearned	Unlearned
1627	Stable Diffusion	by CL [11]	by TARF
1628		the second s	
1629			
1630			
1631			
1632			
1633			
1634			
1635			
1636			
1637			
1638			
1639			
1640			
1641			
1642		Transfer of the second s	Lawy are and
1643			
1644			
1645			
1646			
1647			
1648			- Minn
1649			
1650			
1650		and distre.	Large
1652			And the second
1654			and a second second
1655			
1656		A CORE AL	
1657			
1658		, 5	
1659			
1660			
1661			
1662			
1663			
1664			
1665			
1666		HINE ET CONTRACT	HING ELLONAUX
1667			
1668			
1669			
1670			
1671			
1672			

¹⁶⁷⁴ F ADDITIONAL EXPERIMENTAL RESULTS

- 1676 In this section, we provide additional experimental results of our work.
- 1678 In Appendix F.1, we summarize the additional experimental setups.
- ¹⁶⁷⁹ In Appendix F.2, we discuss the crucial target identification in unlearning.
- In Appendix F.3, we discuss and compare TARF with the advanced method in all matched scenario.
- ¹⁶⁸² In Appendix F.4, we discuss potential ways to extend unlearning to the scenario without class labels.
- In Appendix F.5, we verify unlearning on large-scale datasets trained with large models.
- ¹⁶⁸⁵ In Appendix F.6, we present unlearning with different model structures.
- ¹⁶⁸⁶ In Appendix F.7, we present the full results under multiple runs with the four forgetting tasks.
- 1689 F.1 EXTRA EXPERIMENTAL SETUPS

1690 We introduce additional experimental details in the specific unlearning tasks. In our TARF, In general, 1691 we set $t_1 = 1$ for all the target identification parts, and we adopt k = 0.04, $t_0 = 2$ in model mismatch 1692 forgetting, and k = 0.02, $t_0 = 2$ for all matched, target mismatch and data mismatch forgetting in 1693 the unlearning request on CIFAR-10 classification task; for the CIFAR-100 classification task, we 1694 adopt k = 0.5, $t_0 = 2$ in model mismatch forgetting, and k = 0.05, $t_0 = 2$ for all matched, target 1695 mismatch and data mismatch forgetting. For the other hyperparameters, we follow the previous works [30, 37, 11] to set the specific values. All the forgetting trails use 10 epochs for the total 1697 unlearning process except for GA (use 5 epochs) and IU (use the specific fixed step for optimization). 1698 The specific parameters and the pre-trained models (unlearn base) are provided in our source codes.

1699 1700

1688

F.2 DISCUSSION ABOUT TARGET IDENTIFICATION IN UNLEARNING

1702 In this part, we further discuss the important factors for the achievability of the unlearning tasks. To 1703 be more specific, for the target or data mismatch forgetting, the scenario assumes that the identified 1704 forgetting data is part of the whole samples belonging to the target concept, which means there are 1705 other forgetting data included in the remaining set that need to be found. Thus, target identification is important for effective unlearning. As demonstrated in Section 3.2, the representation gravity 1706 can be a useful clue in forgetting dynamics to identify the other false retaining data. An implicit 1707 assumption is that those false retaining data have similar semantic features to the initially provided 1708 forgetting data, which has smaller representation distance than the retaining part of data as illustrated 1709 in Figure 12. Empirically, the model can have similar prediction changes on those false retaining 1710 data with the initial forgetting data. However, not all of the superclasses officially defined for the 1711 CIFAR-100 dataset are suitable for constructing the unlearning request, as some superclasses are not 1712 semantically separable like "aquatic mammals" and "fish". It can be found in Figure 14, where we 1713 check the Top-10 classes with the most accuracy changes after gradient ascent for each superclass 1714 in the CIFAR-100 dataset, some false retaining data (class-level indicated by blue arrows) are not 1715 easily identified given the two initially provided forgetting data classes (indicated by red arrows). 1716 One interesting future problem can be how to handle the spurious correlation given the insufficient 1717 representative samples.

1718

1719 1720 F.3 DISCUSSION ABOUT TARF ON ALL MATCHED SCENARIO

Regarding the all matched scenario, there is no need for the target identification part to identify extra forgetting data in the all-matched scenario as the target concept matches the forgetting data, then TARF degenerates into a general framework using the given forgetting data to forget, and the rest to retain. The performance of TARF is comparable to the existing best counterpart like SCRUB regarding the "Gap↓" in Table 2. It can be found that the overall performance of the unlearned models has already closely approximated the Retrained reference. Furthermore, since TARF is a general framework, we can also adopt the KL divergence loss with the original model as designed in SCRUB to further improve the performance, for which we present the comparison in Table 17.



1772 1773

1774

F.4 DISCUSSION ABOUT TARF ON WEAKLY-SUPERVISED SCENARIO

Our current work mainly focus on expanding the scope of conventional class-wise unlearning. Regarding the existing approximate unlearning studies [51, 37, 6, 30], considering the all matched forgetting scenario with full supervision, we push it towards more practical settings via decoupling the class labels and the target concept. For machine unlearning under weak supervision, there are limited studies [51] to our best knowledge, and we believe it is worth an in-depth exploration in future work.

Given that if a model is trained with semi-supervised or other weak supervision, we can obtain the pseudo labels by the model prediction for its unlearning phase. Instead of using the predicted label,





CIFAR-100 dataset. Note that some target concepts are not successfully identified by the identified data.

Figure 14: Task Identification using the CIFAR-100 dataset for target mismatch forgetting.

1826 we can also utilize the distillation objective to encourage the unlearned model's output to be far 1827 away from (or close to) the original ones. With the guide of model prediction, the data belonging 1828 to the same superclass with the forgetting data can be figured out to constrain the unlearning target. In Table 18, we present the results of our methods when only the given forgetting data are labeled, 1829 demonstrating our framework can be extended to achieve satisfactory performance. 1830

1831 1832

1834

FORGETTING IN THE LARGE-SCALE DATASET F.5 1833

In this part, we present more experiments conducted on large-scale dataset like ImageNet-1k in 1835 Table 19, and also unlearning multiple classes in the large-scale datasets in Table 20.

-	Type / \mathcal{D}	Dataset		CI	FAR-10					CIF	AR-100		
		Method / Metrics U.	A RA	TA	MIA	Gap↓	TIME↓	UA	RA	TA	MIA	Gap↓	TIME↓
		Retrained (Ref.) 0.0	0 99.51	94.69	100.00	-	43.3	0.00	97.85	76.03	100.00	-	43.2
		FT [58] 1.0	7 98.62 0 00.04	92.36	100.00	1.07	4.43	0.67	96.32	72.34	100.00	1.47	5.02
	Semi-supervised	SCRUB [57] 0.0	0 99.94	91.00	100.00	1.05	2.00	0.00	99.98	70.75	100.00	0.71	3.23
	Scenarios	TARF (with CE)0.0TARF (with KL)0.0	0 98.23 0 98.81	91.95 93.33	100.00 100.00	1.01 0.52	4.21 4.32	$0.00 \\ 0.00$	96.90 96.95	72.53 74.98	100.00 100.00	1.11 0.49	4.68 4.89

1836 Table 17: Performance comparison in the all matched scenario when TARF with CE loss/KL divergence (refer to Eq. (16) with the original model for the retaining part.

Table 18: A case study (%) on the unlearning on CIFAR-100 under the weakly-supervised scenario (e.g., using the pseudo-label generated by model prediction to handle unlabeled retaining data).

Type / D	Dataset	Model mismatch					Data mismatch						
	Method / Metrics	UA	RA	TA	MIA	Gap↓	TIME↓	UA	RA	TA	MIA	Gap↓	TIME↓
	Retrained	88.22	98.58	78.50	25.78	-	43.8	0.00	98.50	80.15	100.00	-	53.2
	FT	92.67	95.02	79.34	16.33	4.58	4.86	82.62	95.66	79.77	37.24	37.15	4.93
	RL	80.11	95.83	79.83	99.00	21.35	4.93	89.78	96.82	79.90	70.76	30.49	4.97
Semi-supervised	GA	6.78	94.83	76.96	97.78	39.68	0.06	6.00	97.65	79.23	98.04	2.43	0.05
Scenarios	BS	18.11	95.90	72.28	95.22	37.14	0.89	15.38	98.50	72.28	96.22	6.76	0.96
	L ₁ -sparse	82.11	85.17	75.22	20.00	7.15	5.00	84.53	85.13	75.22	17.02	46.45	5.03
	TARF (full labels)	86.67	97.05	80.07	26.00	1.21	4.81	0.00	95.01	78.98	100.00	1.17	4.78
	TARF (unlabeled retain)	90.22	96.58	80.01	22.54	2.17	4.84	1.33	95.30	78.12	99.34	1.45	4.85

1855 F.6 FORGETTING WITH DIFFERENT MODEL STRUCTURES

In this part, we further check the unlearning performance of our TARF on different pre-trained model 1857 structures compared with several baselines. We choose CIFAR-100 as the pre-training classification task and conduct all matched forgetting and model mismatch forgetting. The results are summarized 1859 in Table 21. The results validate that our TARF can robustly achieve better unlearning performance across different model structures. 1861

FULL RESULTS WITH DIFFERENT FORGETTING TASKS F 7 1863

1864 In this section, we provide the full results of Table 2, which is conducted by setting different random 1865 seeds (for multiple runs) with the original trails and reported as the mean and std values for each 1866 evaluation metric. Tables 22 to 25 presents the performance of unlearning on CIFAR-10, and Tables 26 1867 to 29 presents the performance of unlearning on CIFAR-100. The performance comparison of our 1868 TARF with other baseline across the four forgetting tasks (i.e., all matched, target mismatch, model 1869 mismatch, and data mismatch) demonstrated the general effectiveness of our algorithm framework. 1870

G **BROADER IMPACT**

1872 1873

1871

1862

1838 1839

1843 1844

1845

1846 1847

1849 1850 1851

In this work, we explore the label domain mismatch in class-wise unlearning, which aims to enhance 1874 the flexibility of data regulation with the increasing concern about the trustworthiness of machine 1875 learning. Pushing forward the practical usage of machine unlearning, our research provides a broader 1876 consideration of real-world unlearning scenarios and offers significant positive social impacts. It can 1877 enhance data privacy protection by allowing individuals to effectively remove their data, ensuring 1878 some sensitive data is not used for analysis. In addition, unlearning can remove bias or discrimination 1879 by correcting flawed datasets, promoting the development of fairness or other ethical considerations. 1880 This feature also enables enterprises to adhere to data protection standards such as GDPR [49] and 1881 CCPA [48], therefore promoting confidence among users. Our newly introduced unlearning setting, which decouples the class label and the target concept, is more general and discusses the achievability 1882 of various unlearning requests, which may often be different from the taxonomy of pre-training tasks. 1883

1884 Although we take a step forward in more practical class-wise unlearning by considering the label 1885 domain mismatch scenarios, it is not the end of this direction and there are still many problems to be addressed. Following the previous works [58, 17, 30, 6], our work mainly focuses on the class-wise unlearning with the classification model for the exploration, future efforts can also be paid in the unlearning problem of the emerging and powerful generative models. On the technical level, although those compared unlearning methods and our framework can achieve the forgetting target, it 1889 all requires extra computational cost, and how to make it more efficient can be further studied.

Type / \mathcal{D}	Dataset			All 1	natched					Target	mismatcl	h	
	Method / Metrics	UA	RA	TA	MIA	Gap↓	TIME↓	UA	RA	TA	MIA	Gap↓	TIME↓
	Retrained	0.00	79.77	77.64	100.00	-	7075.48	0.00	80.09	77.54	100.00	-	7777.54
	FT PL [56]	0.00	70.18	71.98	100.00	3.82	608.11	0.79	70.26	72.07	100.00	4.02	608.62
	GA	0.00	66.25	67.36	19.40	5.95	8.76	0.00	31.21	37.74	0.00	45.14	17.38
	BS	0.00	31.15	36.33	100.00	22.48	9.03	0.00	21.57	27.56	99.97	27.13	23.75
	L_1 -sparse	0.00	67.98	70.70	100.00	4.68	603.21	0.00	67.24	70.28	100.00	5.03	601.27
	SCRUB	29.77	74.92	/3.00	81.//	13.71	033.42	22.44	/4.8/	/3.00	82.77	11./1	081.33
	TARF (ours)	0.00	70.53	72.23	100.00	3.66	600.11	0.00	69.93	71.79	100.00	3.97	628.87
ImageNet-1k	Dataset	Model matched Data mismate								mismatch			
	Method / Metrics	UA	RA	TA	MIA	Gap↓	TIME↓	UA	RA	TA	MIA	Gap↓	TIME↓
	Retrained	79.15	80.00	70.29	25.69	-	6501.27	0.00	80.36	70.38	100.00	-	6493.10
	RL [56]	83.31 87.62	70.38 69.43	64.05 63.26	19.00	6.68 9.13	695.42 959.84	88.21	69.99 70.33	63.81	12.21	4.24 48.15	956.13
	GA [28]	0.00	66.62	58.91	100.00	44.56	17.44	0.00	15.35	14.34	0.00	55.26	17.58
	BS [6]	0.00	45.81	40.84	100.00	54.28	19.69	0.00	13.00	12.10	100.00	31.41	23.70
	L_1 -sparse [30] SCRUB [37]	82.00 86.08	67.94 74.82	62.58 68.04	19.15 14.69	6.34	1091.29 663.61	0.00	66.37 74.84	61.03 67.92	100.00 93.10	5.84 7.27	689.82
	TARF (ours)	80.62	70.27	64.04	19.46	5.92	601.28	0.00	70.10	63.97	100.00	4.17	602.62

Table 19: Results (%). Comparison with the unlearning baselines on ImageNet-1k. All matched forgetting: unlearn 1 class; Target mismatch forgetting: unlearn three classes belonging to "fish".

Table 20: Results (%). Comparison with the unlearning baselines on TinyImageNet-200/ImageNet-1k with more (10+) forgetting classes in all matched forgetting scenarios.

Scenarios / \mathcal{D}	Unlearn request		forget 1	0 classes	in Tiny-l	mageNe	t		forget 3	0 classes	s in Tiny-I	mageNe	t
	Method / Metrics	UA	RA	TA	MIA	Gap↓	TIME↓	UA	RA	TA	MIA	Gap↓	TIME↓
	Retrained FT GA	0.00 2.04 17.76	71.00 70.63 61.74	60.29 59.04 56.12	100.00 98.26 76.90	- 1.35 13.57	251.43 27.10 1.37	0.00 2.79 28.95	65.26 72.41 59.72	57.60 60.36 57.54	100.00 97.38 57.06	- 3.71 19.37	181.13 35.00 3.49
	TARF (ours)	0.00	69.63	59.69	100.00	0.49	28.5	0.00	70.24	60.16	100.00	1.89	39.6
All matched	Unlearn request		forget 5	0 classes	in Tiny-l	mageNe	t		forge	et 10 clas	ses in Ima	ageNet	
Forgetting	Method / Metrics	UA	RA	TA	MIA	Gap↓	TIME↓	UA	RA	TA	MIA	Gap↓	TIME↓
	Retrained FT	0.00 5.19	66.26 75.77	57.88 61.29	100.00 85.75	- 8.09 22.16	161.37 44.62 7.70	0.00	51.94 55.16	56.74 59.53	100.00 100.00	- 1.50	917.66 316.14
	TARF (ours)	0.00	71.68	60.89	100.00	23.10 2.11	46.97	0.00	50.69	55.83	100.00	0.85	353.69

CIFAR-100	Task		1	All matcl	hed		Model mismatch						
011111 100	Metric	UA	RA	TA	MIA	Gap↓	UA	RA	TA	MIA	Gap↓		
	Retrained	0.00	97.26	73.13	100.00	-	87.44	98.22	82.12	19.89	-		
	FT [58]	0.00	90.92	66.86	100.00	3.15	95.22	95.17	77.71	7.56	6.89		
VGG-19	RL [56]	0.00	90.29	66.16	100.00	3.48	96.22	95.26	77.71	98.56	23.7		
	GA [28]	0.00	79.27	56.03	100.00	8.77	0.00	93.09	74.30	100.00	45.13		
	TARF (ours)	0.00	91.96	67.94	100.00	2.62	82.67	93.71	76.24	24.22	4.87		
	Retrained	0.00	97.85	76.03	100.00	-	88.22	98.58	78.50	25.78	-		
	FT [58]	0.66	96.55	71.97	100.00	1.51	98.22	96.79	80.14	6.78	8.11		
ResNet-18	RL [56]	0.11	95.90	71.57	100.00	1.63	94.11	96.70	80.17	96.89	20.14		
	GA [28]	1.89	95.26	69.14	99.89	2.87	9.33	95.13	77.22	96.89	38.68		
	TARF (ours)	0.00	96.90	71.51	100.00	1.37	86.00	96.54	74.20	22.78	2.89		
	Retrained	0.00	97.71	76.95	100.00	-	88.11	98.37	83.61	23.56	-		
	FT [58]	0.67	96.61	71.29	100.00	1.86	97.44	95.70	78.70	7.33	8.29		
WideResNet	RL [56]	0.00	95.86	71.36	100.00	1.86	85.77	94.69	78.26	96.00	20.9		
	GA [28]	0.44	91.49	66.29	100.00	2.26	4.33	91.76	75.18	99.11	43.7		
	TARF (ours)	0.00	96.51	71.77	100.00	1.60	88.00	95.50	79.06	22.67	2.11		

1945	Table 21: Results (%) of unlearning with different model structure. All methods are trained on the
10/6	same backbone, i.e., the basis of unlearning initialization is the same (except for retraining from
1047	scratch). Values are percentages. Bold numbers are superior results. \downarrow indicates smaller are better.
1947	

1966Table 22: Main Results (%). Comparison with the unlearning baselines. All methods are trained on1967the same backbone, i.e., the basis of unlearning initialization is the same (except for retraining from1968scratch). Values are percentages. Bold numbers are superior results. \downarrow indicates smaller are better.

CIFAR-10	Metric	U	4	R	A	T	A	ML	A	Ga	p↓
	Method	mean	std	mean	std	mean	std	mean	std	mean	std
	Retrained	0.00	-	99.51	-	94.69	-	100.00	-	-	-
	FT [58]	4.66	3.59	98.58	0.04	92.42	0.06	100.00	0.00	1.96	0.89
	RL [56]	2.23	1.90	98.30	0.65	91.97	0.74	100.00	0.00	1.54	0.82
A 11 4 - h - J	GA [28]	0.34	0.16	95.48	0.24	88.52	0.35	99.88	0.10	2.67	0.21
	IU [29]	0.11	0.05	72.50	15.65	68.28	14.10	99.98	0.02	13.39	7.41
All matcheu	BS [6]	24.72	0.32	88.91	0.97	81.84	0.94	89.23	0.56	14.74	0.70
	L_1 -sparse [30]	0.00	0.00	94.18	0.03	90.01	0.24	100.00	0.00	2.50	0.05
	SalUn [11]	0.48	0.46	88.66	2.67	84.48	2.40	100.00	0.00	5.39	1.38
	SCRUB [37]	1.23	0.58	99.92	0.02	91.23	0.56	100.00	0.00	1.28	0.23
	TARF (ours)	0.00	0.00	98.22	0.02	92.09	0.14	100.00	0.00	0.97	0.03

Table 23: Main Results (%). Comparison with the unlearning baselines. All methods are trained on the same backbone, i.e., the basis of unlearning initialization is the same (except for retraining from scratch). Values are percentages. Bold numbers are superior results. \downarrow indicates smaller are better.

1986												
1987	CIFAR-10	Metric	U U	4	R4	4	T/	4	MI	A	Gaj	p↓
1988		Method	mean	std								
1989		Retrained	87.76	-	99.58	-	95.91	-	20.57	-	-	-
1990		FT [58]	94.78	0.11	98.65	0.12	93.77	0.21	10.42	0.86	5.06	0.27
1991		RL [56]	48.25	5.43	98.01	0.12	93.03	0.21	98.10	0.64	30.37	1.53
1002		GA [28]	6.49	0.73	86.91	0.08	82.03	0.18	94.39	0.59	45.41	0.27
1992	Model	IU [29]	15.84	7.86	85.89	1.45	81.08	1.49	93.58	3.71	43.36	3.62
1993	mismatch	BS [6]	14.05	3.76	53.28	2.51	51.25	1.86	94.90	1.06	59.75	2.29
1994		L_1 -sparse [30]	92.25	0.87	95.01	0.25	91.67	0.04	17.40	2.86	4.14	1.00
1995		SalUn [11]	16.31	7.40	92.91	1.05	86.50	2.12	99.24	0.09	41.55	2.14
1006		SCRUB [37]	93.21	1.17	99.83	0.13	93.29	0.81	14.24	0.87	3.65	0.18
1990		TARF (ours)	89.91	1.20	97.73	0.24	92.66	0.17	20.36	2.54	2.45	0.46

Table 24: Main Results (%). Comparison with the unlearning baselines. All methods are trained on the same backbone, i.e., the basis of unlearning initialization is the same (except for retraining from scratch). Values are percentages. Bold numbers are superior results. \downarrow indicates smaller are better.

CIFAR-10	Metric	UA		R	A	T.	4	MIA		Gap↓	
011111 10	Method	mean	std	mean	std	mean	std	mean	std	mean	s
	Retrained	0.00	-	99.38	-	93.85	-	100.00	-	-	
	FT [58]	52.23	1.80	98.43	0.05	91.74	0.09	50.59	0.15	26.18	0
	RL [56]	50.63	0.62	98.21	0.65	91.51	0.61	53.88	2.36	25.06	(
	GA [28]	41.64	0.82	97.05	0.04	89.68	0.17	63.23	1.10	21.23	(
Target	IU [29]	45.32	0.81	70.25	17.82	65.67	2.76	55.98	2.76	36.66	9
mismatch	BS [6]	53.78	0.16	89.67	1.02	79.34	3.95	66.31	10.02	25.36	
	L_1 -sparse [30]	49.55	0.08	93.57	0.05	89.06	0.23	51.33	0.09	27.21	(
	SalUn [11]	47.85	1.22	87.84	3.25	83.38	2.94	58.10	2.85	27.40	
	SCRUB [37]	48.53	1.02	99.43	0.21	91.66	0.28	51.27	0.73	24.92	(
	TARF (ours)	0.05	0.02	97.65	0.08	91.28	0.47	100.00	0.00	1.09	(
											_

Table 25: Main Results (%). Comparison with the unlearning baselines. All methods are trained on the same backbone, i.e., the basis of unlearning initialization is the same (except for retraining from scratch). Values are percentages. Bold numbers are superior results. \downarrow indicates smaller are better.

CIFAR-10	Metric	U	4	R.	A	T/	4	MI.	A	Ga	p↓
011111110	Method	mean	std	mean	std	mean	std	mean	std	mean	std
	Retrained	0.00	-	99.53	-	95.56	-	100.00	-	-	-
	FT [58]	96.85	0.06	98.62	0.13	93.47	0.21	6.93	0.45	48.23	0.18
	RL [56]	73.62	2.86	97.90	0.22	92.59	0.66	52.04	2.23	31.55	1.49
	GA [28]	9.82	1.13	96.14	0.28	90.46	0.33	90.46	0.95	6.56	0.67
Data	IU [29]	15.19	7.66	94.80	0.70	89.08	0.46	92.83	4.26	8.39	2.69
mismatch	BS [6]	16.72	0.02	61.01	0.21	53.81	4.05	93.47	1.24	25.88	1.27
	L_1 -sparse [30]	95.42	0.35	94.57	0.26	91.07	0.01	10.82	1.30	48.51	0.47
	SalUn [11]	55.52	3.76	92.68	1.19	89.25	1.22	60.23	3.30	27.12	2.37
	SCRUB [37]	97.06	0.52	99.16	0.23	94.72	0.56	9.98	0.43	46.98	0.21
	TARF (ours)	0.00	0.00	98.35	0.18	93.42	0.34	100.00	0.00	0.83	0.13

Table 26: Main Results (%). Comparison with the unlearning baselines. All methods are trained on the same backbone, i.e., the basis of unlearning initialization is the same (except for retraining from scratch). Values are percentages. Bold numbers are superior results. \downarrow indicates smaller are better.

CIFAR-100	Metric	U.	UA		4	T /	TA		MIA		p↓
	Method	mean	std	mean	std	mean	std	mean	std	mean	sto
	Retrained	0.00	-	97.85	-	76.03	-	100.00	-	-	-
	FT [58]	0.67	0.01	96.44	0.12	72.16	0.19	100.00	0.00	1.49	- 0.0
A 11 4 - h J	RL [56]	0.56	0.45	96.00	0.10	71.79	0.22	100.00	0.00	1.66	0.
	GA [28]	1.61	0.28	95.00	0.26	68.85	0.29	99.89	0.00	2.93	0.
	IU [29]	0.00	0.00	39.80	2.19	31.09	1.51	100.00	0.00	25.75	0.
An matcheu	BS [6]	4.83	0.05	90.17	0.06	64.30	0.64	99.45	0.12	6.20	0.
	L_1 -sparse [30]	0.00	0.00	94.25	0.57	71.35	1.27	100.00	0.00	1.92	0.
	SalUn [11]	0.00	0.00	77.00	1.66	63.06	0.92	100.00	0.00	8.46	0.
	SCRUB [37]	0.00	0.00	99.72	0.26	76.69	0.06	100.00	0.00	0.64	0.
	TARF (ours)	0.00	0.00	96.67	0.24	72.40	0.14	100.00	0.00	1.21	0.

Table 27: Main Results (%). Comparison with the unlearning baselines. All methods are trained on the same backbone, i.e., the basis of unlearning initialization is the same (except for retraining from scratch). Values are percentages. Bold numbers are superior results. \downarrow indicates smaller are better.

CIFAR-100	Metric	U	4	R	4	ТА		MIA		Ga	p
011111 100	Method	mean	std	mean	std	mean	std	mean	std	mean	-
	Retrained	88.22	-	98.58	-	78.50	-	25.78	-	-	Ī
	FT [58]	95.45	2.78	95.91	0.89	79.74	0.40	11.56	4.78	6.34	
	RL [56]	87.11	7.00	96.27	0.44	80.00	0.17	97.95	1.06	20.75	
	GA [28]	8.06	1.28	94.98	0.15	77.09	0.13	97.34	0.45	39.18	
Model	IU [29]	39.95	5.28	97.22	0.39	79.71	0.63	83.28	3.17	27.08	
mismatch	BS [6]	18.56	0.56	95.87	0.03	74.96	2.68	94.95	0.28	36.27	
	L_1 -sparse [30]	91.11	5.00	94.28	0.18	77.61	0.39	15.56	4.45	5.84	
	SalUn [11]	74.78	8.45	79.98	1.14	71.55	0.77	65.61	11.39	19.71	
	SCRUB [37]	92.45	2.80	99.44	0.78	78.75	1.75	20.13	4.56	4.14	
	TAPE (ours)	8178	1.00	07.10	0.14	80.02	0.15	28.80	2.80	2 37	-
		07.70	1.90	71.19	0.14	00.02	0.15	20.09	2.09	4.57	

Table 28: Main Results (%). Comparison with the unlearning baselines. All methods are trained on the same backbone, i.e., the basis of unlearning initialization is the same (except for retraining from scratch). Values are percentages. Bold numbers are superior results. \downarrow indicates smaller are better.

CIFAR-100	Metric	U	4	R	4	T/	4	ML	A	Ga	p↓
	Method	mean	std	mean	std	mean	std	mean	std	mean	std
	Retrained	0.00	-	97.85	-	73.72	-	100.00	-	-	-
	FT [58]	58.58	0.40	96.42	0.10	72.31	0.22	45.94	0.83	28.87	0.34
	RL [56]	57.76	1.14	96.00	0.10	72.04	0.16	50.67	3.69	27.66	1.15
	GA [28]	22.07	0.69	96.87	0.24	70.52	0.30	90.45	0.23	8.95	0.10
Target	IU [29]	30.80	0.18	39.44	2.25	31.00	1.42	63.83	0.14	42.03	0.91
mismatch	BS [6]	40.91	0.47	98.36	0.04	70.04	1.38	85.00	0.16	15.03	0.18
	L_1 -sparse [30]	55.31	2.90	94.23	0.44	72.15	1.27	48.47	3.54	30.26	1.18
	SalUn [11]	43.29	1.60	77.15	1.63	63.30	0.93	64.63	1.34	27.45	0.10
	SCRUB [37]	59.56	0.09	99.74	0.26	76.14	0.82	45.45	0.56	29.60	0.02
	TARF (ours)	0.29	0.03	97.06	0.29	73.27	0.41	100.00	0.00	0.38	0.17

Table 29: Main Results (%). Comparison with the unlearning baselines. All methods are trained on the same backbone, i.e., the basis of unlearning initialization is the same (except for retraining from scratch). Values are percentages. Bold numbers are superior results. \downarrow indicates smaller are better.

CIFAR-100	Metric	UA		R.	A T		4	MIA		Gap↓	
	Method	mean	std	mean	std	mean	std	mean	std	mean	s
	Retrained	0.00	-	98.50	-	80.15	-	100.00	-	-	
	FT [58]	90.79	5.18	96.19	0.52	79.80	0.03	20.46	16.78	43.25	5
	RL [56]	93.60	3.82	96.32	0.39	79.92	0.02	65.20	5.56	32.73	2
	GA [28]	6.98	0.98	97.78	0.14	79.34	0.11	97.53	0.51	2.75	0
Data	IU [29]	37.22	5.71	99.17	0.21	80.01	1.81	85.41	2.41	13.54	2
mismatch	BS [6]	15.71	0.33	98.47	0.04	76.02	3.74	96.05	0.18	5.86	(
	L_1 -sparse [30]	89.02	4.67	94.18	0.05	78.89	0.20	18.67	4.36	41.64	2
	SalUn [11]	79.00	6.07	79.92	1.05	71.55	0.51	44.18	9.96	39.42	3
	SCRUB [37]	93.28	2.10	99.25	0.98	79.18	0.48	18.45	3.55	46.13	2
	TARF (ours)	0.00	0.00	95.80	0.79	79.55	0.57	100.00	0.00	1.61	(
											_