# SHARPNESS-AWARE MACHINE UNLEARNING

**Anonymous authors** 

000

001 002 003

004

006

008 009

010

011

012

013

014

016

017

018

019

021

025

026

027 028

029

031

033

034

035

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

#### **ABSTRACT**

We characterize the effectiveness of Sharpness-aware minimization (SAM) under machine unlearning scheme, where unlearning forget signals interferes with learning retain signals. While previous work prove that SAM improves generalization with noise memorization prevention, we show that SAM abandons such denoising property when fitting the forget set, leading to altered generalization depending on signal strength. We further characterize the signal surplus of SAM in the order of signal strength, which enables learning from less retain signals to maintain model performance and putting more weight on unlearning the forget set. Empirical studies show that SAM outperforms SGD with relaxed requirement for retain signals and can enhance various unlearning methods either as pretrain or unlearn algorithm. Motivated by our refined characterization of SAM unlearning and observing that overfitting can benefit more stringent sample-specific unlearning, we propose Sharp MinMax, which splits the model into two to learn retain signals with SAM and unlearn forget signals with sharpness maximization, achieving best performance. Extensive experiments show that SAM enhances unlearning across varying difficulties measured by memorization, yielding decreased feature entanglement between retain and forget sets, stronger resistance to membership inference attacks, and a flatter loss landscape. Our observations generalize to more noised data, different optimizers, and different architecures.

# 1 Introduction

Deep neural networks have grown so large and complex that retraining a model from scratch to forget even a few samples has become impractically costly in both computation and energy. This challenge has catalyzed the study of machine unlearning: methods that efficiently remove the influence of specific training data without full retraining, aiming to forget designated examples while preserving overall performance. Numerous unlearning strategies have been explored – from influence-based updates that subtract a data point's contribution (Izzo et al., 2021), to fine-tuning with targeted weight sparsification (Jia et al., 2023), to joint optimization approaches that explicitly balance "retain" vs. "forget" objectives by gradient ascent/descent on different data subsets (Kurmanji et al., 2023). However, a fundamental understanding of what makes unlearning effective remains elusive. Key questions persist: How should we trade off forgetting unwanted data versus retaining accuracy on the rest? How do different training algorithms influence unlearning dynamics? Why are some samples inherently harder to forget than others? In practice, the lack of principled answers has led to ad-hoc hyperparameter tuning and unpredictable behavior across tasks. In particular, when a model is simultaneously fed with conflicting retain and forget signals, these signals can interfere and even cancel out during training, hampering the unlearning process (Kurmanji et al., 2023). To date, there are few robust solutions to mitigate this interference, underscoring the need for a deeper theoretical foundation for machine unlearning.

Recent advances in learning theory and optimization hint at possible directions to tackle these issues. First, a signal-versus-noise perspective has provided new insight into model behavior: for example, Chen et al. (2023) formalize how networks learn meaningful patterns while ignoring or memorizing label noise, and Zhao et al. (2024) empirically identify factors that make certain data points harder to forget. Particularly relevant is the Sharpness-Aware Minimization (SAM) method (Foret et al., 2020) that has been shown to seek flatter loss minima and thereby dramatically reduce memorization of noisy data, leading to improved generalization in noisy-label settings (Chen et al., 2023). These observations suggest that a model's ability to distinguish true signal from noise may be key to effective unlearning. An optimizer that naturally suppresses memorization of noise might also

 be better suited for forgetting specific examples when required. To investigate this hypothesis, we quantify each sample's memorization level using established metrics (Feldman, 2020; Feldman & Zhang, 2020), allowing us to rank the "forget set" by difficulty. This enables a controlled study of how different optimization algorithms perform when asked to forget data that the model has learned to varying extents.

We present a comprehensive theoretical and empirical study of machine unlearning through the combined lens of signal-noise decomposition and sharpness-aware optimization. We focus on the challenging scenario where both retain and forget samples are present in each training batch with mixed objectives, and we compare standard Stochastic Gradient Descent (SGD) to SAM in this context. Building on recent theoretical frameworks for ReLU networks (Kou et al., 2023), we derive rigorous results for a two-layer CNN that characterize the unlearning process under each optimizer. Our analysis yields several striking findings. (1) SAM's noise suppression can break down under unlearning: we prove that when tasked with intentionally forgetting a set of samples (treated as "noise"), SAM is forced by objective to abandon its usual denoising behavior – effectively overfitting to the forget set nearly as much as SGD does. This result challenges the expectation that flatter-minima methods would inherently excel at unlearning. (2) We establish formal guidelines for balancing retain vs. forget objectives: in particular, we derive the minimum retain-weighting factor  $\alpha$  needed to prevent catastrophic forgetting of the kept data. Our theory shows that SAM can accomplish successful unlearning with a significantly smaller retain weight  $\alpha$  than SGD, meaning SAM tolerates a stronger forgetting signal without sacrificing retained accuracy. In the regime of benign overfitting (where the model fits even noisy data without large generalization error), we quantify the gap in required  $\alpha$  between SAM and SGD and prove it scales on the order of  $O(\sqrt{d/n})$  (with d the model dimension and n the training set size). (3) Perhaps most surprisingly, our findings call for a re-examination of overfitting in unlearning. Contrary to conventional wisdom, we show that deliberate overfitting – in a controlled way that limits its impact on the rest of the data – can enhance the complete removal of those samples. This insight is especially relevant in stringent privacy or copyright scenarios, suggesting that the strict avoidance of overfitting may not always be optimal.

Our contributions can be summarized as follows:

**Theoretical Framework:** We introduce a rigorous analytical framework for machine unlearning based on signal-noise decomposition. This framework explicitly models the interplay between retain and forget signals. Using this lens, we analyze the behaviors of SGD versus SAM and prove that SAM's denoising advantage "shuts off" on forget data: when SAM is asked to unlearn labeled noise, it ends up overfitting to the forget set almost as much as SGD.

Balancing Retain vs. Forget Objectives: We derive provable guidelines for balancing the retain/forget trade-off. In particular, we identify the minimal value of the weighting ratio parameter  $\alpha$  that guarantees sufficient retention of knowledge. We show that SAM requires a strictly smaller  $\alpha$  than SGD to achieve effective unlearning. In the regime of benign overfitting for both the optimizers, we analytically bound the difference in required  $\alpha$  on the order of  $O(\sqrt{d/n})$ .

**Empirical Validation:** Through extensive experiments on CIFAR-100 and ImageNet datasets, we validate our theoretical insights. We demonstrate that incorporating SAM into state-of-the-art unlearning methods consistently boosts forgetting efficacy while better preserving accuracy on the remaining data. Models optimized with SAM yield flatter loss landscapes and reduced entanglement between retained and forgotten samples, corroborating our theory that SAM distinguishes signal from noise better. We also observe that SAM-trained models are less vulnerable to membership inference attacks to forget set, indicating improved unlearning.

**Novel Unlearning Algorithm:** Finally, inspired by our analysis, we propose Sharp MinMax, a new unlearning approach that decouples the retain and forget objectives. Sharp MinMax splits the model into two cooperative parts: one is trained with SAM on the retained data, while the other performs sharpness maximization on the forget data to intentionally overfit those samples to ensure forgottenness. This design mitigates interference between retain and forget signals. Sharp MinMax achieves state-of-the-art unlearning performance in our experiments, especially on challenging high-memorization forget sets, where it significantly outperforms existing techniques in completely erasing the target data's influence.

# 2 PRELIMINARIES

#### 2.1 Data and Model Construction

We construct a practical learning scenario which distinguishes between useful and unrelated signals from inputs. Similar constructions have been adopted in previous work (Kou et al., 2023; Chen et al., 2023) with rich notation. For convenience, we summarize a table of notation in App. C. Consider learning binary classification with label  $y \in \{\pm 1\}$  using a two-layer CNN on image training data set  $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i \in [n]} \sim \mathcal{D}$ . Each image consists of P patches and assign randomly one of them as the signal  $y_i \varphi$  for label  $y_i$  and the universal signal vector  $\varphi \in \mathbb{R}^d$ , and represent other patches by the noise vector  $\xi_i \in \mathbb{R}^d \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$ . Thus, each input image is vectorized as  $\mathbf{x}_i = [\xi_i, ..., y_i \varphi, ..., \xi_i] \in \mathbb{R}^{P \times d}$ , where  $y_i \varphi$  can appear at any position.

The second layer of CNN is fixed as  $\pm 1/m$  respectively for m convolutional filters. The two-classes network can be expressed as  $f(\mathbf{W}, \mathbf{x}) = f_{+1}(\mathbf{W}_{+1}, \mathbf{x}) - f_{-1}(\mathbf{W}_{-1}, \mathbf{x})$ , where

$$f_j(\mathbf{W}_j, \mathbf{x}) = \frac{1}{m} \sum_{r=1}^m \sum_{p=1}^P \sigma(\langle \mathbf{w}_{j,r}, \mathbf{x} \rangle) = \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{j,r}, y\varphi \rangle) + (P-1)\sigma(\langle \mathbf{w}_{j,r}, \xi \rangle). \tag{1}$$

Here  $\sigma$  denotes ReLU activation,  $\mathbf{w}_{j,r} \in \mathbb{R}^d$  denotes the weight for the r-th filter, and  $\mathbf{W}_j$  is the collection of model weights for  $j=\pm 1$ . We train this CNN with cross-entropy loss  $\mathcal{L}(\mathbf{W},\mathcal{S})$ . Denote  $\mathbf{w}_{j,r}^{(t,b)}$  for  $j \in \{\pm 1\}, r \in [m]$  the convolutional filter at the b-batch of t-th epoch of SGD.

We decompose the weight update into learning signal and noise coefficients  $\kappa_{j,r}^{(t,b)}$ ,  $\zeta_{j,r,i}^{(t,b)}$  for learning the signal and the noise respectively, such that

$$\mathbf{w}_{j,r}^{(t,b)} = \mathbf{w}_{j,r}^{(0,0)} + j \cdot \kappa_{j,r}^{(t,b)} \cdot \boldsymbol{\varphi} \| \boldsymbol{\varphi} \|_{2}^{-2} + (P-1)^{-1} \sum_{i=1}^{n} \zeta_{j,r,i}^{(t,b)} \cdot \boldsymbol{\xi}_{i} \| \boldsymbol{\xi}_{i} \|_{2}^{-2},$$
(2)

where the learning goal is to increase  $\kappa_{j,r}^{(t,b)}$  and decrease  $\zeta_{j,r,i}^{(t,b)}$ . This construction also extends to multiclass classification considering one vs. all setting with K binary classification problems.

## 2.2 SIGNAL-TO-NOISE UNLEARNING

Given a pretrained model  $f_{\mathcal{A}}^{T_1}$  by algorithm  $\mathcal{A}$  for  $T_1$  epochs on  $\mathcal{S}$ , machine unlearning aims to eliminate the influence of forget set  $\mathcal{F} \subseteq \mathcal{S}$  to the model training, while maintain generalizability to unseen data without compromising performance on the remaining retain set  $\mathcal{R} = \mathcal{S} \setminus \mathcal{F}$ . Denote the unlearned model as  $f_{\mathcal{A}}^{T_2}$  by unlearning algorithm  $\mathcal{U}$ , which is initialized as  $f_{\mathcal{A}}^{T_1}$  and unlearned for  $T_2$  epochs. We consider unlearning a small portion of  $\mathcal{S}$  with much less expense than retraining the model from scratch on  $\mathcal{R}$ , so  $|\mathcal{F}| < |\mathcal{R}|$  and  $T_2 < T_1$ .

**Random Label** (RL) (Graves et al., 2021) aims to unlearn by finetuning on  $\mathcal{S}$  but with  $\mathcal{F}$ 's labels randomly flipped in each epoch. It naturally fits into our setup as label-flipped  $\mathcal{F}$  become the noise, and motivates us to investigate unlearning algorithms under the same theoretical framework. The gradient update of  $\kappa_{i,r}^{(t,b)}$  and  $\zeta_{i,r,i}^{(t,b)}$  can be expressed as

$$\kappa_{j,r}^{(t,b+1)} = \kappa_{j,r}^{(t,b)} - \frac{\eta \|\boldsymbol{\varphi}\|_{2}^{2}}{Bm} \left[ \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{R}}} \ell_{i}^{\prime(t,b)} \sigma'(\langle \mathbf{w}_{j,r}^{(t,b)}, \widehat{y}_{i} \boldsymbol{\varphi} \rangle) - \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{F}}} \ell_{i}^{\prime(t,b)} \sigma'(\langle \mathbf{w}_{j,r}^{(t,b)}, \widehat{y}_{i} \boldsymbol{\varphi} \rangle) \right],$$

$$\zeta_{j,r,i}^{(t,b+1)} = \zeta_{j,r,i}^{(t,b)} - \frac{\eta (P-1)^{2} \|\boldsymbol{\xi}_{i}\|_{2}^{2}}{Bm} \cdot \ell_{i}^{\prime(t,b)} \sigma'(\langle \mathbf{w}_{j,r}^{(t,b)}, \boldsymbol{\xi}_{i} \rangle) \cdot \operatorname{sgn}(y_{i} = j),$$
(3)

where  $B, \eta$  denote the batch size and learning rate,  $\operatorname{sgn}(\cdot)$  denotes  $\pm 1$  sign function,  $\mathcal{I}^{\mathcal{R}}_{t,b}$  and  $\mathcal{I}^{\mathcal{F}}_{t,b}$  denote batch samples from  $\mathcal{R}$  and  $\mathcal{F}$ , respectively. In each iteration,  $\mathcal{I}^{\mathcal{F}}_{t,b}$  aims to erase its signal in  $\kappa^{(t,b)}_{j,r,i}$ , while  $\xi_i$  reinforces or decreases  $\zeta^{(t,b)}_{j,r,i}$  update depending on label agreement.

**Negative Gradient** (NegGrad) (Kurmanji et al., 2023) actively unlearns  $\mathcal{F}$  using gradient ascent while gradient-descending on  $\mathcal{R}$ . The combined loss objective is defined as

$$\mathcal{L}_{\text{NegGrad}}(\mathbf{W}, \mathcal{R}, \mathcal{F}) = \frac{1}{|\mathcal{R}|} \sum_{i \in \mathcal{R}} \alpha \ell \left( y_i f\left(\mathbf{W}, \mathbf{x}_i\right) \right) - \frac{1}{|\mathcal{F}|} \sum_{i \in \mathcal{F}} (1 - \alpha) \ell \left( y_i f\left(\mathbf{W}, \mathbf{x}_i\right) \right). \tag{4}$$

Minimizing  $\mathcal{L}_{\text{NegGrad}}$  induces competing gradients, canceling each other during  $\kappa, \zeta$  update.  $\alpha$  serves as a weight coefficient that accounts for the size imbalance between  $\mathcal{R}$  and  $\mathcal{F}$ . To synchronously optimize the model with retain and forget samples, we draw B samples from both subsets each batch and train for  $|\mathcal{R}|/B$  batches. Thus, forget samples' signals are relatively enlarged by a fraction of  $|\mathcal{R}|/|\mathcal{F}|$  due to repetition. Heuristically,  $\alpha \propto |\mathcal{R}|/(|\mathcal{F}|+|\mathcal{R}|)$ .

#### 2.3 Denoising Property of SAM

Sharpness-Aware Minimization (SAM) (Foret et al., 2020) aims to minimize a perturbed empirical loss at the worst point in the neighborhood of **W**, solving the following optimization problem:

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{W}, \mathcal{S}) + \left[ \max_{\widehat{\epsilon}} \mathcal{L}(\mathbf{W} + \widehat{\epsilon}, \mathcal{S}) - \mathcal{L}(\mathbf{W}, \mathcal{S}) \right], \tag{5}$$

for a controlled perturbation  $\hat{\epsilon}$ . It ensures a uniformly low training loss and avoids sharp landscape. While both SGD and SAM learn a sufficient signal with  $\kappa_{j,r}^{T_1} = \Omega(1)$  after  $T_1$  epochs, Chen et al. (2023) prove that SAM outperforms SGD by noise suppression and SAM upper bounds  $\zeta_{j,r,i}^{T_1}$  by O(1) while SGD is dimension dependent  $O(\log d)$ . The key difference stems from the noise memorization prevention of SAM. Given the perturbation term  $\hat{\epsilon}^{(t,b)}$  in SAM:

$$\widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)} = \frac{\tau}{m} \sum_{i \in \mathcal{I}_{t,b}} \sum_{p \in [P]} \ell_i^{\prime(t,b)} j \cdot y_i \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t,b)}, \mathbf{x}_{i,p} \rangle) \mathbf{x}_{i,p} \cdot \left\| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}^{(t,b)}, \mathcal{I}_{t,b}) \right\|_F^{-1}, \tag{6}$$

consider ReLU activation at any fixed iterate  $\mathbf{w}_{j,r}^{(t,b)}$  for SGD:  $\langle \mathbf{w}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \rangle \geq 0$  vs. SAM:  $\langle \mathbf{w}_{j,r}^{(t,b)} + \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \rangle$  for  $k \in \mathcal{I}_{t,b}, j = y_k$ . SAM's  $\langle \mathbf{w}_{j,r}^{(t,b)} + \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \rangle$  expands to  $\langle \mathbf{w}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \rangle + \langle \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \rangle$ , where  $\langle \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \rangle$  is proven to be sufficiently negative to cancel  $\langle \mathbf{w}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \rangle$  by selecting a proper  $\tau$ , thus deactivating the noise (Chen et al., 2023). This effectively prevents SAM from learning from the noise which would lead to harmful overfitting for SGD. We are curious about whether SAM improves unlearning: a flatter landscape can make learning easier, then it should make unlearning easier too despite a reverse sign. But is it a simple adaptation, and can we straightforwardly extend previous theories and findings to develop unlearning algorithms?

## 3 SHARPNESS-AWARE UNLEARNING

We first show that the SAM's noise memorization prevention in Sec. 2.3 does not fully hold when SAM is used with NegGrad for gradient ascent on  $\mathcal{F}$ . Specifically, SAM overfits to forget signals as much as SGD, while maintaining its denoising property on  $\mathcal{R}$ . Based on this result, we derive refined test error bounds for SGD and SAM under NegGrad and characterize the different  $\alpha$  thresholding between SGD and SAM for unlearning. Although SAM continues to improve unlearning and maintain generalizability, the altered activation patterns and unlearning behaviors are not captured by previous works, as SAM is forced to fit forget signals (viewed as noise) by NegGrad objective. This leads to divergent behaviors on  $\mathcal{R}$  and  $\mathcal{F}$ , which can be of independent interest.

## 3.1 NegGrad Revisited

Unlike RL, the mutual interference between  $\mathcal{F}$  and  $\mathcal{R}$  under NegGrad additionally affects  $\zeta_{j,r}$  update. The update rules for  $\kappa_{j,r}^{(t,b)}$  and  $\zeta_{j,r}^{(t,b)}$  under NegGrad now become:

$$\kappa_{j,r}^{(t,b+1)} = \kappa_{j,r}^{(t,b)} - \frac{\eta \|\boldsymbol{\varphi}\|_{2}^{2}}{Bm} \left[ \alpha \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{R}}} \nabla_{\boldsymbol{\varphi}_{i}} - (1-\alpha) \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{F}}} \nabla_{\boldsymbol{\varphi}_{i}} \right],$$

$$\zeta_{j,r}^{(t,b+1)} = \zeta_{j,r}^{(t,b)} - \frac{\eta(P-1)^{2}}{Bm} \left[ \alpha \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{R}}} \nabla_{\boldsymbol{\xi}_{i}} - (1-\alpha) \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{F}}} \nabla_{\boldsymbol{\xi}_{i}} \right],$$
(7)

where  $\nabla_{\varphi_i} = \ell_i'^{(t,b)} \sigma'(\langle \mathbf{w}_{j,r}^{(t,b)} + \delta, y_i \varphi \rangle), \nabla_{\xi_i} = \operatorname{sgn}(y_i = j) \|\xi_i\|_2^2 \ell_i'^{(t,b)} \sigma'(\langle \mathbf{w}_{j,r}^{(t,b)} + \delta, \xi_i \rangle),$  and  $\delta = \hat{\epsilon}_{j,r}^{(t,b)}$  for SAM and 0 for SGD. In plain words, a sample  $i \in \mathcal{R}$  of class j causes a decrease in  $\zeta_{j,r,i}$ , discouraging memorizing noise for the correct class, while another sample  $i' \in \mathcal{R}$  of class -j

causes an increase in  $\zeta_{j,r,i}$ , encouraging  $w_{j,r}$  to use  $\xi_i$  to distinguish class j from -j. Conversely, a sample  $i \in \mathcal{F}$  of class j, which we want to predict -j in unlearning, will increase  $\zeta_{j,r,i}$ , encouraging  $w_{j,r}$  to use noise  $\xi_i$  in a way that harms class j, and vice versa. Similar intuition also applies to  $\kappa_{j,r}$ . The interference in  $\zeta_{j,r}$  update will alter SAM's behaviors towards forget signals as summarized in Lemma 3.1.

**Lemma 3.1** (Noise memorization of  $\mathcal{F}$  by SAM under NegGrad). Under the NegGrad scheme and the Assumption D.1 holds, we have that if for SGD:  $\langle \mathbf{w}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \rangle \geq 0, k \in \mathcal{I}_{t,b}^{\mathcal{R}}$  and  $j = y_k$ , then for SAM:  $\langle \mathbf{w}_{j,r}^{(t,b)} + \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \rangle < 0$ . However, if for SGD:  $\langle \mathbf{w}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \rangle \geq 0, k \in \mathcal{I}_{t,b}^{\mathcal{F}}$  and  $j = y_k$ , then for SAM:  $\langle \mathbf{w}_{j,r}^{(t,b)} + \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \rangle > 0$ .

See proof in App. D.2. Because the activation patterns on  $\mathcal{I}_{t,b}^{\mathcal{R}}$  and  $\mathcal{I}_{t,b}^{\mathcal{F}}$  diverge, SAM continues to suppress noise memorization and leverage its sharpness-aware updates when fitting  $\mathcal{R}$ , but "falls back" to SGD-like behavior on  $\mathcal{F}$ . This split yields two distinct sets of bounds on  $\kappa_{j,r}$  and  $\zeta_{j,r}$  for  $\mathcal{R}$  and  $\mathcal{F}$ , which lead to separate test errors shown in App. D.1 and D.2. However, given a pretrained model  $f_{\mathcal{A}}^{T_1}$  with  $\kappa_{j,r}^{T_1} > 0$  to start unlearning, as long as retain signals weighted by  $\alpha$  dominate, the signal strength will remain sufficient and continue to grow. This is shown in Chen et al. (2023) when the signal strength is saturated at  $T < T_1$ . We can thus choose  $\alpha$  threshold based on this principle. With proper forget-retain size ratio, results in Chen et al. (2023) still hold: SGD's test error converges when signal strength is sufficient, but can't be upper bounded otherwise; SAM's test error converges either way.  $\beta$  serves as a knob to control the convergence rate:

**Theorem 3.2** (SGD test error under NegGrad). Under Assumption D.1, for any  $\epsilon > 0$  and  $1 > \alpha \ge |\mathcal{R}|/(|\mathcal{F}| + |\mathcal{R}|) := \beta > 0.5$ , then with probability at least  $1 - \delta$ , the training loss converges:  $\mathcal{L}(\mathbf{W}^T, \mathcal{D}) \le \epsilon$ . Moreover, if  $\|\varphi\|_2 \ge C_1 d^{1/4} n^{-1/4} P \sigma_p$ , we have the test error  $\mathcal{L}^{test}(\mathbf{W}^T, \mathcal{D}) \le \epsilon$ . If  $\|\varphi\|_2 \le C_3 d^{1/4} n^{-1/4} P \sigma_p$ , we have  $\lim_{\beta \to 0.5} \mathcal{L}^{test}(\mathbf{W}^{T_2}, \mathcal{D}) \ge 0.1$ , and  $\lim_{\beta \to 0.5} \mathcal{L}^{test}(\mathbf{W}^{T_2}, \mathcal{D}) \ge 0.05$ .

**Theorem 3.3** (SAM test error under NegGrad). Under Assumption D.I, for any  $\epsilon > 0$  and  $1 > \alpha \ge |\mathcal{R}|/(|\mathcal{F}| + |\mathcal{R}|) := \beta > 0.5$ , choose  $\tau = \Theta(\frac{m\sqrt{B}}{P\sigma_p\sqrt{d}})$ . Then with probability at least  $1 - \delta$ , the training loss converges:  $\mathcal{L}(\mathbf{W}^T, \mathcal{D}) \le \epsilon$ . Moreover, if  $\|\varphi\|_2 \ge C_1 d^{1/4} n^{-1/4} P\sigma_p$ , we have  $\lim_{\beta \to 1} \mathcal{L}^{test}(\mathbf{W}^T, \mathcal{D}) \le \epsilon$ . If  $\Omega(1) \le \|\varphi\|_2 \le C_3 d^{1/4} n^{-1/4} P\sigma_p$ : we still have  $\lim_{\beta \to 1} \mathcal{L}^{test}(\mathbf{W}^T, \mathcal{D}) \le \epsilon$ .

See proofs in App. D.1 and D.2. Together, these theorems describe how SGD and SAM behave when retain signals dominate. For SAM, if  $\|\varphi\|_2 \leq C_3 d^{1/4} n^{-1/4} P \sigma_p$ , it will suffer harmful overfitting to  $\mathcal{F}$ . However, as long as  $\alpha \geq |\mathcal{R}|/(|\mathcal{F}| + |\mathcal{R}|)$  and  $\|\varphi\|_2 \geq \Omega(1)$ , learning on  $\mathcal{R}$  guarantees overall benign training and yields a bounded test error. Under the same condition, Corollary 3.3.1 concludes that while the signal coefficient continues to grow for both SGD and SAM, SGD's noise accumulation is loosely bounded by model dimension, while SAM's by O(1):

**Corollary 3.3.1**  $(\kappa, \zeta \text{ update under NegGrad})$ . Under the NegGrad, if  $\alpha \geq |\mathcal{R}|/(|\mathcal{F}| + |\mathcal{R}|)$ , since  $\kappa_{j,r}^{T_1} = \Omega(1)$ , both SGD and SAM continue to grow. Given the learned  $\zeta_{j,r}^{T_1}$ , SGD continues to overfit the noise with  $O(\log d)$ , while SAM overfit the noise from  $\mathcal{F}$  with  $O(\log d)$  and from  $\mathcal{R}$  with O(1).

See proof in App. D.3. Finally, we characterize the differed choice of  $\alpha$  for SGD and SAM as SAM learns signal more efficiently. We also reveal that  $\alpha$  depends not only on forget-retain size ratio as commonly conjectured, but also on the signal strength, and thus the dimensionality of the problem:

**Lemma 3.4** (Signal-surplus of SAM under NegGrad). Under the NegGrad, for any  $\varphi$  where  $\|\varphi\|_2 \geq \Omega(1)$ , SAM exhibits faster signal learning on  $\mathcal{R}$ :  $\Delta^{SAM}_{epoch}\kappa_{j,r}/\Delta^{SGD}_{epoch}\kappa_{j,r} = \Theta(\|\varphi\|_2^2)$ .

See proof in App. D.4. As a result, SAM relies on a more relaxed  $\alpha$  threshold than SGD due to faster signal learning. For SGD to achieve the same signal learning performance as SAM, we need to scale up  $\alpha^{\text{SGD}}$  to satisfy  $\alpha^{\text{SGD}}/\alpha^{\text{SAM}} = \Theta(\|\varphi\|_2^2)$ . If  $\|\varphi\|_2 \ge C_1 d^{1/4} n^{-1/4} P \sigma_p$  and both SGD and SAM achieve benign overfitting, then given the extra signal learning from  $\mathcal{R}$ , SAM results in faster  $\kappa$  update and a surplus signal of  $\Theta(d^{1/2}|\mathcal{R}|^{-1/2}P^2\sigma_p^2)$  in each unlearning epoch.

## 3.2 SHARP MINMAX

In Sec. 3.1, we showed that SAM is provably better on out of sample test errors under NegGrad, and we empirically verify that SAM achieves better unlearning performance in Sec. 4. But how does the refined characterization matter, given maintained test error conclusions? Jointly with empirical observations, the altered behaviors of SAM on  $\mathcal F$  motivates new unlearning algorithms. Our experiments show that SAM+NegGrad attains higher forget accuracy than SGD+NegGrad, forgetting less effectively. This finding forces us to reconsider the conventional view that overfitting is always detrimental: while overfitting indeed harms generalization, it may be beneficial when the goal is to remove specific samples from a model. Consequently, for abstract concept forgetting we continue to demand strong generalization; but for stringent scenarios—where exact sample removal is mandated by privacy or compliance constraints—a model's tendency to overfit can actually enhance its unlearning of those exact points. The divergent behaviors under SAM+NegGrad motivates the following new algorithm: we can split a portion of model parameters to purposefully overfit to  $\mathcal{F}$ , while leaving the rest to maximally maintain the model utility by leveraging SAM purely on  $\mathcal{R}$ . Motivated by how SGD with sharper minima tends to forget better, we propose Sharp MinMax to intentionally optimize for sharper-than-SGD minima with the purpose of overfitting to forget signals for unlearning. Inspired by Kim et al. (2023), we leverage sharpness maximization:

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{W}, \mathcal{S}) - \left[ \max_{\widehat{\epsilon}} \mathcal{L}(\mathbf{W} + \widehat{\epsilon}, \mathcal{S}) - \mathcal{L}(\mathbf{W}, \mathcal{S}) \right], \tag{8}$$

resulting in a sharper landscape that harms the generalization by overfitting. We then apply weight masking based on gradient magnitudes to divide our model into two during optimization. Specifically, we pass  $\mathcal{F}$  to  $f_{\mathcal{A}}$  once, accumulate gradients for each parameter, and check top parameters with smallest magnitudes cut off by a given percentage. Smaller gradient magnitudes suggest more fitting to the forget samples during the pretraining stage, which demands more unlearning. We then apply SAM on the retain model and sharpness maximization on the forget model. The retain model with SAM is already characterized by Chen et al. (2023), while the forget model requires a stronger signal strength than SGD to avoid harmful overfitting. See implementation details in App. E.2.

## 3.3 QUANTIFYING UNLEARNING DIFFICULTY WITH MEMORIZATION

We examine the effectiveness of unlearning  $\mathcal{U}$  based on memorization, which sufficiently reveals the difficulty of unlearning (Zhao et al., 2024). Feldman & Zhang (2020) define the degree to which a sample is memorized by a pretraining  $\mathcal{A}$  on example ( $\mathbf{x}_i, y_i$ ) from  $\mathcal{S}$  as the memorization score:

$$\operatorname{mem}(\mathcal{A}, \mathcal{S}, i) := \Pr_{f \leftarrow \mathcal{A}(\mathcal{S})} \left[ f\left(\mathbf{W}, \mathbf{x}_i\right) = y_i \right] - \Pr_{f \leftarrow \mathcal{A}(\mathcal{S} \setminus i)} \left[ f\left(\mathbf{W}, \mathbf{x}_i\right) = y_i \right], \tag{9}$$

where  $\mathcal{S}\setminus i$  denotes  $\mathcal{S}$  with the sample  $(\mathbf{x}_i,y_i)$  removed. Samples of high-memorization scores can be atypical samples which model usually learns later in the training process after more updates to the model than typical ones. Thus unlearning them would be harder and may require more iterations of unlearning steps which may impact the model performance on the retain distribution. The converse is true for samples of low-memorization scores. We can hence construct  $\mathcal{F}$  of varying unlearning difficulties based on memorization scores to comprehensively evaluate  $\mathcal{U}$ .

## 4 EMPIRICAL STUDY

We conduct major experiments on CIFAR-100 (Krizhevsky et al., 2009) and ImageNet-1K (Russakovsky et al., 2015) using ResNet-50 (He et al., 2016), and adopt pre-computed memorization scores for from Feldman & Zhang (2020) to generate  $\mathcal F$  of different difficulties with  $|\mathcal F| \approx 5\% |\mathcal S|$ , denoted as  $[\mathcal F_{high}, \mathcal F_{mid}, \mathcal F_{low}]$ . For both pretraining and unlearning, we adopt SAM (Foret et al., 2020) with  $\rho=0.1$  and Adaptive SAM (ASAM) (Kwon et al., 2021) with  $\rho=[0.1,1.0]$ . We ensure same optimal hyper-paprameters for each comparable [SGD,SAM] pair. See details in App. E.

**Evaluation.** We follow previous work (Triantafillou et al., 2024; Zhao et al., 2024) to measure the tug-of-war tradeoff between forgetting and retaining of  $f_{\mathcal{U}}$  based on accuracy  $\mathrm{Acc}(\theta, \mathcal{D})$ , with the retrained model  $f_{\mathcal{A}(\mathcal{R})}$  as reference:

$$ToW(f_{\mathcal{U}}) = (1 - (Acc(f_{\mathcal{A}(\mathcal{R})}, \mathcal{R}) - Acc(f_{\mathcal{U}}, \mathcal{R}))) \cdot (1 - (Acc(f_{\mathcal{U}}, \mathcal{F}) - Acc(f_{\mathcal{A}(\mathcal{R})}, \mathcal{F})))$$

$$\cdot (1 - (Acc(f_{\mathcal{A}(\mathcal{R})}, \mathcal{D}_{test}) - Acc(f_{\mathcal{U}}, \mathcal{D}_{test}))), \text{ with test transforms on } \mathcal{R}, \mathcal{F}.$$

$$(10)$$

Table 1:  $\operatorname{ToW}(\%) \uparrow$  of unlearning on ImageNet-1K and CIFAR-100. For each  $(\mathcal{U}, \mathcal{A})$  pair, we report ToW of each  $\mathcal{F}$  and compute averages. SAM consistently improves current unlearning methods.

ImageNet			SGD			A = AS	AM 0.1				AM 1.0				AM 0.1	
Unlearn $\mathcal{U}$	High	Mid	Low	AVG												
NegGrad	78.764	84.199	88.515	83.826	78.426	83.93	86.651	83.002	78.522	83.929	89.947	84.133	78.03	84.176	88.839	83.682
+ASAM 0.1	78.52	84.113	89.188	83.94	78.366	84.07	89.098	83.845	78.762	84.267	90.579	84.536	78.083	84.062	89.973	84.039
+ASAM 1.0	78.966	83.389	92.174	84.843	78.975	83.358	91.843	84.725	78.027	83.326	92.772	84.708	77.762	83.284	92.617	84.554
+SAM 0.1	77.898	82.985	92.841	84.575	78.301	83.04	91.722	84.354	77.388	82.473	93.429	84.43	76.807	82.587	92.829	84.074
RL	74.598	86.617	86.714	82.643	74.857	86.462	86.192	82.504	74.317	86.813	87.630	82.92	74.055	86.715	88.594	83.121
+ASAM 1.0	74.951	85.581	91.069	83.867	75.221	85.473	90.425	83.707	73.950	85.393	91.516	83.62	73.579	85.494	91.74	83.604
SalUn	44.981	71.839	95.008	70.609	46.104	71.735	94.652	70.83	45.814	72.308	95.116	71.079	46.006	72.419	95.218	71.214
+ASAM 1.0	45.998	71.554	95.628	71.06	46.938	71.268	95.224	71.143	45.856	71.695	95.924	71.158	46.358	72.034	95.791	71.394
CIFAR100			SGD			A = AS					AM 1.0				AM 0.1	
Unlearn $\mathcal{U}$	High	Mid	Low	AVG												
NegGrad	78.334	83.335	83.718	81.796	79.277	86.454	88.637	84.789	77.274	78.59	85.443	80.436	67.826	74.145	76.374	72.78
+ASAM 0.1	78.131	82.846	86.78	82.586	80.336	87.539	87.671	85.182	77.331	79.074	88.039	81.482	70.054	74.158	78.087	74.1
+ASAM 1.0	80.806	81.465	87.052	83.108	82.196	84.391	90.502	85.696	78.731	79.264	93.249	83.748	72.518	75.653	86.759	78.31
+SAM 0.1	81.331	75.059	94.151	83.514	82.86	77.94	94.179	84.993	74.704	70.898	95.898	80.5	65.080	66.089	95.078	75.416
L1-Sparse	63.448	68.686	53.991	62.042	63.699	72.775	60.34	65.605	61.252	68.197	61.47	63.64	65.258	71.941	59.014	65.404
+ASAM 1.0	66.903	75.554	58.967	67.141	66.213	77.119	66.697	70.01	65.117	73.754	62.517	67.129	63.051	74.556	65.117	67.575
SCRUB	58.418	76.125	12.708	49.084	67.163	79.09	10.823	52.359	57.816	73.176	58.483	63.158	43.246	68.433	17.368	43.016
+ASAM 1.0	50.313	73.353	97.631	73.766	60.515	80.204	97.508	79.409	48.569	73.09	97.776	73.145	18.137	61.618	97.933	59.229
RL	68.464	84.395	72.4	75.086	64.518	80.215	69.711	71.481	66.689	86.411	69.677	74.259	64.391	85.481	70.55	73.474
+ASAM 1.0	69.952	86.779	74.409	77.047	66.909	86.557	69.375	74.280	69.73	91.124	80.321	80.392	72.884	88.633	78.066	79.861
SalUn	69.926	83.056	71.73	74.904	66.541	83.377	71.95	73.956	67.355	89.768	79.095	78.739	69.671	90.495	75.281	78.482
+ASAM 1.0	73.268	92.225	88.175	84.556	71.426	89.182	86.13	82.246	67.715	93.401	89.289	83.468	70.933	92.914	86.477	83.441

Table 2: MIA (%)  $\downarrow$  correctness to  $\mathcal{F}$  on CIFAR-100. We enhance each  $\mathcal{U}$  with ASAM 1.0 and observe consistent improvement.

		A =	SCD		ı	A = AS	AM 0.1		1	1 - 15	AM 1.0			A = SA	M 0 1	
Unlearn U	High	Mid	Low	AVG	High	Mid	Low	AVG	High	Mid	Low	AVG	High	Mid	Low	AVG
L1-Sparse	94.733	63.233	8.6	55.522	94.933	61.367	4.0	53.433	93.833	62.067	5.8	53.9	92.867	60.033	5.033	52.644
+ASAM 1.0	94.267	58.5	5.5	<b>52.756</b>	94.3	57.3	3.6	<b>51.733</b>	93.633	56.033	3.9	<b>51.189</b>	93.8	59.333	3.8	<b>52.311</b>
SCRUB	55.433	18.6	32.6	35.544	64.733	23.1	71.633	53.155	54.767	16.133	9.833	26.911	39.3	9.833	56.3	35.144
+ASAM 1.0	46.467	14.867	0.1	<b>20.478</b>	57.367	22.633	0.167	<b>26.722</b>	44.7	14.567	0.2	<b>19.822</b>	14.433	2.333	0.2	<b>5.655</b>
RL	90.767	62.933	10.767	54.822	91.633	68.267	13.5	57.8	89.067	63.567	15.8	56.145	89.167	61.967	8.267	53.134
+ASAM 1.0	90.3	61.3	9.467	<b>53.689</b>	91.6	62.667	12.7	<b>55.656</b>	88.0	61.3	10.667	<b>53.322</b>	86.3	59.833	5.833	<b>50.655</b>
SalUn	83.433	59.233	7.333	50.0	84.533	59.1	11.167	51.6	79.3	54.667	8.8	47.589	81.467	53.133	6.867	47.156
+ASAM 1.0	79.1	51.833	4.5	<b>45.144</b>	81.7	54.167	6.633	<b>47.50</b>	74.967	49.5	4.2	<b>42.889</b>	75.633	47.667	4.067	<b>42.456</b>
NegGrad	86.933	37.233	2.167	42.111	88.867	40.2	1.733	43.60	82.167	32.1	1.8	38.689	74.667	36.967	3.433	38.356
+ASAM 1.0	84.5	30.1	0.733	<b>38.444</b>	85.6	30.1	0.7	<b>38.8</b>	81.233	24.533	0.533	<b>35.433</b>	73.967	20.733	0.366	<b>31.689</b>

Thus, we encourage high retain/test accuracies and low forget accuracy. Note that our ToW differs from that in previous work as we measure the raw accuracy difference instead of the absolute difference, because new unlearning methods that continue to fine-tune on  $\mathcal R$  can outperform  $f_{\mathcal A(\mathcal R)}$  within a conventional unlearning time  $T_2$ . If using the absolute ToW, a higher test accuracy than  $f_{\mathcal A(\mathcal R)}$  will be penalized and the model performance cannot be properly measured.

## 4.1 SAM CONSISTENTLY OUTPERFORMS WITH BETTER TRADEOFF

We conduct unlearning with various unlearning algorithms  $\mathcal U$  given different pretrained  $f_{\mathcal A}$ . Tab. 1 reports ToW scores of  $\mathcal U$  on CIFAR-100 and ImageNet. We observe that SAM consistently improves all unlearning methods under different initializations  $f_{\mathcal A}^{T_1}$ , suggesting that SAM can universally enhance prevailing  $\mathcal U$ . While different  $\mathcal U$  exhibit varied effectiveness to  $[\mathcal F_{\text{high}}, \mathcal F_{\text{mid}}, \mathcal F_{\text{low}}]$ , we observe that NegGrad achieves a better balance between three forget sets than other methods. We include detailed [retain, forget, test] accuracies, further analysis and demonstration of statistical significance in App. F. Upon close examination on those accuraices, we observe that despite SAM outperforms SGD by better retain and test accuracies and thus better ToW, SGD can oftentimes achieve lower forget accuracies. This aligns with our theoretical analysis where SGD overfits more to  $\mathcal F$ , and it also sparks our Sharp MinMax. Smaller experiments on CIFAR-10 and Tiny-ImageNet in App. G yield aligned conclusions.

MIA correctness. We report correctness rates of membership inference attack (MIA) to  $\mathcal{F}$  on CIFAR-100 in Tab. 2. Lower correctness means better unlearning: forget samples behave more like samples that were never in  $\mathcal{S}$ . We find that SAM consistently improves data privacy while unlearning more effectively. Note that NegGrad achieves better MIA correctness than RL; this is because gradient ascent actively erases gradient signatures of  $\mathcal{F}$  in the model. SCRUB (Kurmanji et al., 2023) with SAM achieves best MIA performance.

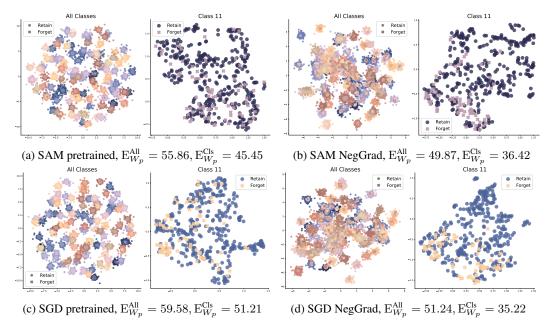


Figure 1: UMAP feature visualization of  $\mathcal{F}_{high}$  on CIFAR-100. We visualize inter-classes and intraclass movements, and class 11 is the largest class in  $\mathcal{F}_{high}$ . For all classes,  $\mathcal{F}$  are assigned to wrong class clusters after NegGrad unlearning. For class-wsie, forget samples gather more tightly.

Our observations further generalize. We consider structured noise unlearning, where another source of noise is introduced during unlearning. We adopt the glass blur and snow effect from ImageNet-C (Hendrycks & Dietterich, 2019) to corrupt  $\mathcal{R}$  and  $\mathcal{F}$  of CIFAR-100, and unlearn with NegGrad and Sharp MinMax. We record experiment results in App. G.3, and observe consistent conclusions where SAM outperforms under both corruptions. We also experiment on ViT-Small (Dosovitskiy et al., 2020) with AdamW (Loshchilov & Hutter, 2017) on CIFAR-100 in App. G.4, with NegGrad and Sharp MinMax. We continue to observe promising improvement by adding SAM, with significant increase of ToW on Sharp MinMax.

#### 4.2 Constrained Overfitting Benefits Unlearning

Table 3:  $ToW(\%) \uparrow$  of Sharp MinMax on ImageNet-1K and CIFAR-100. Comparing with Tab. 1, Sharp MinMax achieves new best ToW performance.

ImageNet		A =	SGD			A = AS	AM 0.1			A = AS	AM 1.0		A =SAM 0.1  High Mid Low AVG				
Unlearn $\mathcal{U}$	High	Mid	Low	AVG	High	Mid	Low	AVG	High	Mid	Low	AVG	High				
SGD	73.357	80.881	86.334	80.191	73.418	80.784	84.378	79.527	73.103	81.105	86.402	80.204	73.052	80.913	85.517	79.827	
ASAM 0.1	78.066	87.914	87.338	84.44	79.077	87.4	86.953	84.476	70.148	88.039	87.554	81.914	78.529	87.642	86.668	84.28	
ASAM 1.0	86.658	87.345	89.694	87.899	86.166	87.192	89.138	87.498	86.915	87.27	90.142	88.109	86.272	87.076	90.064	87.804	
SAM 0.1	86.463	86.755	90.005	87.741	85.511	86.635	89.852	87.333	86.849	86.722	91.111	88.227	85.712	86.486	90.207	87.468	
CIFAR100		A =	SGD			A = A	SAM 0.1			A = AS	SAM 1.0		A = SAM 0.1				
Unlearn $\mathcal{U}$	High	Mid	Low	AVG	High	Mid	Low	AVG	High	Mid	Low	AVG	High	Mid	Low	AVG	
SGD	70.7668	76.692	82.853	76.771	72.137	77.864	81.847	77.282	65.925	74.526	80.127	73.526	60.478	71.931	73.843	68.751	
ASAM 0.1	78.895	96.027	83.473	86.132	84.968	96.451	82.883	88.101	81.825	93.786	87.151	87.587	72.897	80.104	86.659	79.887	
ASAM 1.0	82.27	94.913	86.504	87.896	77.576	99.422	85.894	87.631	84.521	87.761	84.381	85.554	76.037	83.633	77.461	79.044	
SAM 0.1	90.578	90.960	92.494	91.344	91.695	95.543	91.508	92.915	88.664	88.646	93.163	90.158	85.195	78.286	90.963	84.814	

We present ToW of Sharp MinMax and compare to Tab. 1. Compared with NegGrad and other methods, Sharp MinMax further improves the unlearning capabilities across all settings by a noticeable margin, especially on  $\mathcal{F}_{high}$ , and SAM 0.1 achieves ToW > 0.9 for most settings on CIFAR-100. The effectiveness of Sharp MinMax assures our assumptions about overfitting for sample-specific unlearning, providing new insights for designing future unlearning algorithms. By constraining overfitting to only a small portion of model parameters which are most salient to  $\mathcal{F}$ , Sharp MinMax effectively boosts unlearning performance.

## 4.3 QUANTITATIVE ANALYSIS AND VISUALIZATIONS

**Measuring entanglement.** We measure the entanglement between  $\mathcal{R}$  and  $\mathcal{F}$  before and after unlearning. At a coarse level, we implement variance-based entanglement from Goldblum et al. (2020);

Table 4: Entanglement  $\downarrow$  between  $\mathcal{F}$  and  $\mathcal{R}$  of different memorization levels given models based on SGD and ASAM 1.0. While  $E_{Var}$  is hard to conclude a comparison between SGD and SAM across different  $\mathcal{U}$ , SAM shows less entanglement both before and after unlearning than SGD by  $E_{W_n}$ .

SGD		Varian	ce E <sub>Var</sub>			Wasserst	tein $E_{W_p}$		SAM		Varian	ce E <sub>Var</sub>			Wassers	tein $E_{W_p}$	
Model	High	Mid	Low	AVG	High	Mid	Low	AVG	Model	High	Mid	Low	AVG	High	Mid	Low	AVG
Pretrained	30.5	95.28	32.39	52.72	59.58	66.3	63.13	63.0	Pretrained	29.56	88.43	28.91	48.97	55.86	61.74	59.84	59.15
-per class	2.5	6.71	2.51	3.91	51.21	57.11	59.64	55.99	-per class	2.88	6.66	2.71	4.08	45.45	49.88	52.46	49.26
NegGrad	18.87	37.16	22.12	26.05	51.24	52.99	56.12	53.45	NegGrad	17.78	37.49	24.47	26.58	49.87	52.36	54.93	52.39
-per class	0.56	1.8	2.69	1.68	35.22	46.91	55.93	46.02	-per class	0.66	2.03	2.88	1.86	36.42	44.71	50.83	43.99
MinMax	17.7	38.03	21.51	25.75	51.12	53.7	56.77	53.86	MinMax	16.35	32.07	20.75	23.06	51.26	51.8	55.08	52.71
-per class	0.69	2.41	2.27	1.79	38.41	49.57	57.15	48.38	-per class	0.49	1.52	2.97	1.66	33.65	44.56	52.55	43.59

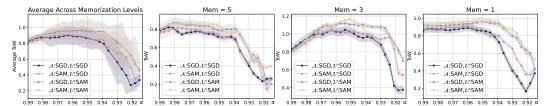


Figure 2: As  $\alpha$  decreases, NegGrad puts less weight on retain signals and learns more from  $\mathcal{F}$ , leading to harmful overfitting. SAM exhibits more tolerance to insufficient retain signals, while  $\mathcal{A}, \mathcal{U} = \text{SGD}$  collapses the fastest. Note that ToW starts failing before  $\alpha = |\mathcal{R}|/(|\mathcal{F}| + |\mathcal{R}|)$ , implying more factors affecting  $\alpha$  threshold as we point out.

Zhao et al. (2024):  $E_{\text{Var}}^{\text{All}}(\mathcal{R}, \mathcal{F}, f) = (\frac{1}{|\mathcal{R}|} \sum_{i \in \mathcal{R}} (\phi_i - \mu_{\mathcal{R}})^2 + \frac{1}{|\mathcal{F}|} \sum_{j \in \mathcal{F}} (\phi_j - \mu_{\mathcal{F}})^2)/((\mu_{\mathcal{R}} - \mu)^2 + (\mu_{\mathcal{F}} - \mu)^2))$ , where  $\phi_i$ ,  $\phi_j$  denote sample embedding,  $\mu_{\mathcal{R}}$ ,  $\mu_{\mathcal{F}}$  denote mean embedding of  $\mathcal{R}$ ,  $\mathcal{F}$ , and  $\mu$  denotes mean embedding over  $\mathcal{R} \cup \mathcal{F}$ . We also compute the class-wise entanglement and report weighted averaged  $E_{\text{Var}}^{\text{Cls}}$ . However,  $E_{\text{Var}}$  assumes good/convex shapes of clusters and relies heavily on cluster means. Inspired by Optimal Transport literature, we propose a refined geometry-aware entanglement based on Wasserstein distance to measure the separation of retain and forget features,  $E_{W_p}^{\text{All}}$  and  $E_{W_p}^{\text{Cls}}$ , which computes the cost of transferring one shaped distribution to another point-wisely. From Tab. 4, we observe that both SGD and SAM unlearning have decreased entanglement with  $E^{\text{Cls}} < E^{\text{All}}$ . While  $E_{\text{Var}}$  cannot further differentiate, we observe that SAM achieves better  $E_{W_p}$  than SGD at all levels. Fig. 1 visualizes the feature space of  $\mathcal{A}$ ,  $\mathcal{U} = \text{ASAM } 1.0$  and  $\mathcal{A}$ ,  $\mathcal{U} = \text{SGD}$  on  $\mathcal{F}_{\text{high}}$ . For all classes, we observe forget samples are assigned to wrong class clusters after NegGrad. For class-wise, we visualize the largest class in  $\mathcal{F}_{\text{high}}$  and observe forget samples to gather more tightly. See App. H.2 for complete visualizations.

**Reducing retain signal.** We verify Lemma 3.4 by reducing  $\alpha$  in NegGrad. Fig. 2 shows ToW changes as  $\alpha$  decreases for various  $\mathcal{A}, \mathcal{U}$  pairs at different memorization levels on CIFAR-100. We observe that  $\mathcal{A}, \mathcal{U} = \text{SGD}$  fails the fastest and hardest, while  $\mathcal{A}, \mathcal{U} = \text{ASAM 1.0}$  exhibits the best resilience. Also note that for CIFAR-100,  $|\mathcal{R}|/(|\mathcal{F}| + |\mathcal{R}|) \approx 0.93$ , but unlearning starts to fail at a higher  $\alpha$ . This verifies our claim that  $\alpha$  depends more than retain-forget ratio.

**Loss landscape.** We visualize loss landscapes of SGD and ASAM 1.0 by perturbing original model along two directions with filter normalization (Li et al., 2018). While SAM unlearning generally keeps flatter landscapes, we observe intriguing phenomena which indicate that unlearning might be an implicit regularizer. See full visualizations and more details in App. H.1.

## 5 Conclusion

In this paper, we provide a refined characterization of SAM under NegGrad unlearning, and theoretical insights on bounding and choosing the weight factor to balance retain and forget signals. Extensive studies verify our analysis and reveals more underlying properties of SAM that are desired for unlearning. Based on our rethinking of overfitting, we also propose a new algorithm which further pushes the boundary of sample-specific unlearning. Our theoretical and empirical findings shed light on future design of unlearning algorithms.

## REFERENCES

- Idan Attias, Gintare Karolina Dziugaite, Mahdi Haghifam, Roi Livni, and Daniel M Roy. Information complexity of stochastic convex optimization: Applications to generalization and memorization. *arXiv preprint arXiv:2402.09327*, 2024.
- Peter L Bartlett, Philip M Long, and Olivier Bousquet. The dynamics of sharpness-aware minimization: Bouncing across ravines and drifting towards wide minima. *arXiv* preprint arXiv:2210.01513, 2022.
- Stella Biderman, Usvsn Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raff. Emergent and predictable memorization in large language models. *Advances in Neural Information Processing Systems*, 36:28072–28090, 2023.
- Nicholas Carlini, Ulfar Erlingsson, and Nicolas Papernot. Distribution density, tails, and outliers in machine learning: Metrics and applications. *arXiv preprint arXiv:1910.13427*, 2019.
- Zixiang Chen, Junkai Zhang, Yiwen Kou, Xiangning Chen, Cho-Jui Hsieh, and Quanquan Gu. Why does sharpness-aware minimization generalize better than sgd? *Advances in neural information processing systems*, 36:72325–72376, 2023.
- Eli Chien, Haoyu Wang, Ziang Chen, and Pan Li. Langevin unlearning: A new perspective of noisy gradient descent for machine unlearning. *arXiv preprint arXiv:2401.10371*, 2024.
- Luc Devroye, Abbas Mehrabian, and Tommy Reddad. The total variation distance between high-dimensional gaussians with the same mean. *arXiv preprint arXiv:1810.08693*, 2018.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. *arXiv* preprint arXiv:2310.12508, 2023.
- Chongyu Fan, Jiancheng Liu, Alfred Hero, and Sijia Liu. Challenging forgets: Unveiling the worst-case forget sets in machine unlearning. In *European Conference on Computer Vision*, pp. 278–297. Springer, 2024.
- Chongyu Fan, Jinghan Jia, Yihua Zhang, Anil Ramakrishna, Mingyi Hong, and Sijia Liu. Towards llm unlearning resilient to relearning attacks: A sharpness-aware minimization perspective and beyond. *arXiv preprint arXiv:2502.05374*, 2025.
- Vitaly Feldman. Does learning require memorization? a short tale about a long tail. In *Proceedings* of the 52nd Annual ACM SIGACT Symposium on Theory of Computing, pp. 954–959, 2020.
- Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. *Advances in Neural Information Processing Systems*, 33:2881–2891, 2020.
- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.
- 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*, pp. 9304–9312, 2020.
- Micah Goldblum, Steven Reich, Liam Fowl, Renkun Ni, Valeriia Cherepanova, and Tom Goldstein. Unraveling meta-learning: Understanding feature representations for few-shot tasks. In *International Conference on Machine Learning*, pp. 3607–3616. PMLR, 2020.
  - Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac machine learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11516–11524, 2021.

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv* preprint arXiv:1903.12261, 2019.
- 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*, pp. 2008–2016. PMLR, 2021.
- Jinghan Jia, Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, Pranay Sharma, and Sijia Liu. Model sparsity can simplify machine unlearning. *Advances in Neural Information Processing Systems*, 36:51584–51605, 2023.
- Pham Khanh, Hoang-Chau Luong, Boris Mordukhovich, and Dat Tran. Fundamental convergence analysis of sharpness-aware minimization. *Advances in Neural Information Processing Systems*, 37:13149–13182, 2024.
- Young In Kim, Pratiksha Agrawal, Johannes O Royset, and Rajiv Khanna. On memorization and privacy risks of sharpness aware minimization. *arXiv preprint arXiv:2310.00488*, 2023.
- Myeongseob Ko, Henry Li, Zhun Wang, Jonathan Patsenker, Jiachen Tianhao Wang, Qinbin Li, Ming Jin, Dawn Song, and Ruoxi Jia. Boosting alignment for post-unlearning text-to-image generative models. *Advances in Neural Information Processing Systems*, 37:85131–85154, 2024.
- Yiwen Kou, Zixiang Chen, Yuanzhou Chen, and Quanquan Gu. Benign overfitting in two-layer relu convolutional neural networks. In *International conference on machine learning*, pp. 17615–17659. PMLR, 2023.
- Alex Krizhevsky et al. Learning multiple layers of features from tiny images. https://www.cs.toronto.edu/kriz/learning-features-2009-TR.pdf, 2009.
- Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. Towards unbounded machine unlearning. *Advances in neural information processing systems*, 36:1957–1987, 2023.
- Jungmin Kwon, Jeongseop Kim, Hyunseo Park, and In Kwon Choi. Asam: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks. In *International Conference on Machine Learning*, pp. 5905–5914. PMLR, 2021.
- Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss land-scape of neural nets. *Advances in neural information processing systems*, 31, 2018.
- Hao Li, Di Huang, Ziyu Wang, and Amir M Rahmani. Skewed memorization in large language models: Quantification and decomposition. *arXiv preprint arXiv:2502.01187*, 2025.
- Yong Liu, Siqi Mai, Xiangning Chen, Cho-Jui Hsieh, and Yang You. Towards efficient and scalable sharpness-aware minimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12360–12370, 2022.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- Peng Mi, Li Shen, Tianhe Ren, Yiyi Zhou, Xiaoshuai Sun, Rongrong Ji, and Dacheng Tao. Make sharpness-aware minimization stronger: A sparsified perturbation approach. *Advances in Neural Information Processing Systems*, 35:30950–30962, 2022.
- USVSN Sai Prashanth, Alvin Deng, Kyle O'Brien, Jyothir SV, Mohammad Aflah Khan, Jaydeep Borkar, Christopher A Choquette-Choo, Jacob Ray Fuehne, Stella Biderman, Tracy Ke, et al. Recite, reconstruct, recollect: Memorization in lms as a multifaceted phenomenon. *arXiv preprint arXiv:2406.17746*, 2024.

- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
  - Yan Scholten, Stephan Günnemann, and Leo Schwinn. A probabilistic perspective on unlearning and alignment for large language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=51WraMid8K.
- Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. Remember what you want to forget: Algorithms for machine unlearning. *Advances in Neural Information Processing Systems*, 34:18075–18086, 2021.
- Eleni Triantafillou, Peter Kairouz, Fabian Pedregosa, Jamie Hayes, Meghdad Kurmanji, Kairan Zhao, Vincent Dumoulin, Julio Jacques Junior, Ioannis Mitliagkas, Jun Wan, et al. Are we making progress in unlearning? findings from the first neurips unlearning competition. *arXiv* preprint *arXiv*:2406.09073, 2024.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Sheng-Yu Wang, Aaron Hertzmann, Alexei Efros, Jun-Yan Zhu, and Richard Zhang. Data attribution for text-to-image models by unlearning synthesized images. *Advances in Neural Information Processing Systems*, 37:4235–4266, 2024.
- Alexander Warnecke, Lukas Pirch, Christian Wressnegger, and Konrad Rieck. Machine unlearning of features and labels. *arXiv preprint arXiv:2108.11577*, 2021.
- Lei Wu, Zhanxing Zhu, et al. Towards understanding generalization of deep learning: Perspective of loss landscapes. *arXiv preprint arXiv:1706.10239*, 2017.
- Yimeng Zhang, Xin Chen, Jinghan Jia, Yihua Zhang, Chongyu Fan, Jiancheng Liu, Mingyi Hong, Ke Ding, and Sijia Liu. Defensive unlearning with adversarial training for robust concept erasure in diffusion models. *Advances in Neural Information Processing Systems*, 37:36748–36776, 2024.
- Zhe Zhang and Guanghui Lan. Solving convex smooth function constrained optimization is almost as easy as unconstrained optimization. *arXiv preprint arXiv:2210.05807*, 2022.
- Kairan Zhao, Meghdad Kurmanji, George-Octavian Bărbulescu, Eleni Triantafillou, and Peter Triantafillou. What makes unlearning hard and what to do about it. *Advances in Neural Information Processing Systems*, 37:12293–12333, 2024.

648 649	A	Rela	ated Works	14
650		A.1	Machine Unlearning	14
651		A.2	Sharpness Aware Minimization	14
652 653		A.3	Data Memorization	14
654 655	В	State	ements	15
656		B.1	Reproducibility Statement	15
657 658		B.2	LLM Usage Statement	15
659		B.3	Limitations and Future Work	15
660 661	C	Tabl	le of Notations	16
662 663				
664	D	Deta	niled Formulations and Proofs	17
665		D.1	Proof to Theorem 3.2	18
666 667		D.2	Proof to Theorem 3.3	24
668			D.2.1 Proof to Lemma 3.1	24
669		D.3	Proof to Corollary 3.3.1	28
670 671		D.4	Proof to Lemma 3.4	29
672				
673	E	Imp	lementation Details	29
674 675		E.1	Experiment Setup	29
676		E.2	Sharp MinMax Implementation	30
677		E.3	Unlearning Setup for Previous Work	30
678 679		E.4	Evaluation Details	31
680 681	F	Deta	niled Empirical Results	31
682		F.1	Statistical Significance	31
683 684		F.2	Complete Accuracies	32
685	~			•
686 687	G		itional Experiments	33
688		<b>G</b> .1	CIFAR-10	33
689		G.2	Tiny-ImageNet	33
690		G.3	Unlearning with Structured Noise	34
691 692		G.4	SAM with Adam and ViT	34
693		~		
694	Н		nplete Visualizations	35
695 696		H.1	Loss Landscape	35
697		H.2	Feature Visualization	36

## A RELATED WORKS

#### A.1 MACHINE UNLEARNING

A wide variety of unlearning algorithms have been proposed to erase the influence of specific data in the pre-trained model. Basic approaches involve finetuning on retain set to unlearn the forget samples with catastrophic forgetting, randomly labeling forget set to force the model to ignore the noisy forget samples, and explicitly "learning to unlearn" from the forget set via gradient ascent (Golatkar et al., 2020; Graves et al., 2021; Warnecke et al., 2021). Recent work pushes the boundaries of each genre with more advanced tools. L1-Sparse (Jia et al., 2023) finetunes on retain set with L1 penalty to improve unlearning with sparsification, NegGrad and SCRUB (Kurmanji et al., 2023) combines gradient descent on retain set and gradient ascent on forget set to jointly update the model, Influence Unlearning and Saliency Unlearning (Izzo et al., 2021; Fan et al., 2023) aim to find model parameters which are important to the forget set for more effective unlearning while preserving model performance. Theoretical work in unlearning draws insights from differential privacy and characterizes distributional closeness in  $(\epsilon, \delta)$ -language. Sekhari et al. (2021) studies unlearning with secondorder update which computes Hessian inverse. Langevin Unlearning (Chien et al., 2024) studies approximate unlearning with privacy and efficiency guarantees based on projected noisy gradient descent. Unlearning also extends to generative vision and language tasks, addressing privacy and safety concerns, erasing concepts, and aligning with human preference (Ko et al., 2024; Wang et al., 2024; Zhang et al., 2024; Scholten et al., 2025).

## A.2 SHARPNESS AWARE MINIMIZATION

Sharpness-aware minimization (SAM) perturbs the model within a ball neighborhood to maximize the loss. Since perturbations in sharp regions result in higher penalties, SAM learns to avoid sharp landscapes and improve generalization with flatness. Recent work improves SAM's flexibility and efficiency. Adaptive SAM (Kwon et al., 2021) introduces scale-invariant adaptive sharpness to address parameter re-scaling sensitivity. GA-SAM (Zhang & Lan, 2022) adapts the perturbation based on gradient strength to improve generalization performance. Sparse SAM (Mi et al., 2022) shows that adding sparsity in perturbations can preserve or even improve performance while accelerating training. LookSAM (Liu et al., 2022) efficiently scales up SAM by only periodically computing the inner gradient ascent. Theoretical studies of SAM focus both on the convergence analysis (Khanh et al., 2024) and its dynamics (Bartlett et al., 2022). Chen et al. (2023) reveal the fundamental mechanism of SAM that prevents memorizing noisy signals by deactivating neurons based on a practical signal-to-noise analytical framework. This inspires us to investigate the intriguing properties of SAM in machine unlearning, where signals from the forget set can be naturally modeled as the noise from the perspective of maintaining model performance with remaining samples.

## A.3 DATA MEMORIZATION

Recent work aims to identify key factors that affect the difficulty of an unlearning task. Fan et al. (2024) define and seek the "worst-case" forget set using a gradient-based adversarial approach. Carlini et al. (2019) investigates and quantifies the atypical-ness of data samples under a differential privacy setting. Zhao et al. (2024) discovers that the more memorized the forget examples are, the harder unlearning becomes. We agree with the empirical studies in Zhao et al. (2024) and study the unlearning effectiveness under different levels of data memorization. Memorization literature provides fundamental understanding and interpretation of learning dynamics and model behaviors, characterizing generalization bounds and the interplay with data (Feldman & Zhang, 2020; Attias et al., 2024). Recent studies also investigate the effects of memorization in large-scale scenarios such as language models (Biderman et al., 2023; Prashanth et al., 2024; Li et al., 2025). Specifically, the memorization and influence scores in Feldman (2020); Feldman & Zhang (2020) provide insights into evaluating unlearning algorithms and designing new approaches. In our study, we have observed varied effectiveness of each unlearning method with respect to forget sets of different memorization levels, and aim at designing unlearning methods which perform well on forgets sets of all difficulties.

## B STATEMENTS

#### **B.1** REPRODUCIBILITY STATEMENT

**Experiment environment.** Our code is built upon several open-source code bases <sup>1</sup> and will be released. We perform all experiments on single NVIDIA A100/H100. We fix random seed for all data processing, saved precomputation (e.g., indices for data subsetting, weight masks), model splitting, pretraining and retraining for reproducible observations. For unlearning parameters and settings, we run experiments with multiple seeds to evaluate statistical significance, see App. F.1.

**Theoretical Assumptions.** Our theoretical analysis follows standard, existing assumptions of model size, data size, effective information in the data (signal) and Gaussian noise in data, which were previously stated in Kou et al. (2023); Chen et al. (2023). In addition to mentioned common assumptions, our Assumption D.1 also assumes conventional unlearning schemes: cross-entropy loss, ReLU activation, clean labels and reasonable size of forget set (< 1/2 trainset size).

## B.2 LLM USAGE STATEMENT

We use GPT to fix grammar and polish short phrases to sharpen our expression. We also use GPT as a smart search engine to gather recent work of interest and summarize existing bug fixes. Zero LLM usage for any core component of our work, including data processing, implementation and experiment, theory, etc., and LLM does not guide the development of any module. No "vibe coding" and mathematical derivation from LLM.

## **B.3** LIMITATIONS AND FUTURE WORK

There are a few limitations based on the signal-to-noise framework, which on the other hand inspire us for future studies. First, there are more interference which can be modeled as noise in machine unlearning, such as the overlap between retain set and forget set. Using hard-cutoff or random sampling to build  $\mathcal F$  might split two similar samples into two opposite subsets, causing interference and impacting unlearning effectiveness. We hypothesize that less overlap between  $\mathcal R$  and  $\mathcal F$  results in more effective unlearning, and vice versa. With more identified and modeled noise sources, another limitation comes from the uncharacterized behaviors when retain signals are weak for some upper bound. Will SAM fail into harmful overfitting under this circumstance? Theoretical and empirical studies under this situation might leverage the interplay between all signals, including different noisy signals. From an empirical perspective, further analysis of the interactions between and model splitting ratio for Sharp MinMax can be developed, as both factors control the impact of retain and forget signals. Last, we observe an intriguing "regularizing" effect of unlearning using SGD via loss landscape visualization, which demands deeper investigation in future work.

https://github.com/kairanzhao/RUM, https://github.com/davda54/sam, https://github.com/OPTML-Group/Unlearn-Saliency, https://pluskid.github. io/influence-memorization/

## C TABLE OF NOTATIONS

Table 5:

Symbol	Meaning / Notes	Symbol	Meaning / Notes
$\mathbf{x}_i \in \mathbb{R}^{P \times d}$	Input image of sample $i$ , vectorized into $P$ patches of dimension $d$ (one patch holds the signal $y_i \varphi$ and $P-1$ patches contain noise)	$y_i \in \{\pm 1\}$	Binary class label for sample $i$
$\boldsymbol{\varphi} \in \mathbb{R}^d$	Universal signal vector shared across samples	$\boldsymbol{\xi}_i \in \mathbb{R}^d$	Noise vector for sample $i$ , often drawn from $\mathcal{N}(0, \sigma_p^2 \mathbf{I})$
P	Number of patches per input image	d	Dimensionality of each patch and each convolutional filter
m	Number of convolutional filters per class	$\mathbf{w}_{j,r} \in \mathbb{R}^d$	Weight vector for the $r$ -th filter of class $j \in \{\pm 1\}$
$\mathbf{W}_{j}$	Collection of filters $\{\mathbf w_{j,r}\}_{r=1}^m$ for class $j$	$\mathbf{W}$	Complete set of model parameters
$f(\mathbf{W}, \mathbf{x})$	Two-class CNN output: $f_{+1}(\mathbf{W}_{+1}, \mathbf{x}) - f_{-1}(\mathbf{W}_{-1}, \mathbf{x})$	$f_j(\mathbf{W}_j, \mathbf{x})$	Class-j output: $\frac{1}{m} \sum_{r=1}^{m} \sum_{p=1}^{P} \sigma(\langle \mathbf{w}_{j,r}, \mathbf{x}_{p} \rangle)$
$\sigma(\cdot)$	ReLU activation function	$\sigma'(\cdot)$	Derivative of ReLU used in gradients
$\mathcal{L}(\mathbf{W},\mathcal{S})$	Cross-entropy loss over training set ${\cal S}$	$\ell_i^{\prime(t,b)}$	Gradient of the loss for sample $i$ a epoch $t$ , batch $b$
$\mathbf{w}_{j,r}^{(t,b)}$	$r\text{-th}$ filter of class $j\in\{\pm 1\}$ after $t$ epochs and $b$ batches	$\kappa_{j,r}^{(t,b)}$	Learned signal coefficient for filte $(j, r)$ at step $(t, b)$
$\zeta_{j,r,i}^{(t,b)}$	Learned noise coefficient from sample $i$ on filter $(j, r)$ at step $(t, b)$	$\ oldsymbol{arphi}\ _2, \ oldsymbol{\xi}_i\ _2$	Euclidean norms of the signal and noise vectors
$\mathcal{F}\subseteq\mathcal{S}$	Forget set whose influence is to be removed	$\mathcal{R} = \mathcal{S} \setminus \mathcal{F}$	Retain set used for continued training
$f_{\mathcal{A}}^{T_1}$	Model after $T_1$ epochs of training by algorithm $\mathcal{A}$	$f_{\mathcal{U}}^{T_2}$	Model after $T_2$ epochs of unlearning by algorithm $\mathcal{U}$
$T_1, T_2$	Numbers of epochs for pretraining and unlearning	$\mathcal{I}^{\mathcal{R}}_{t,b},\mathcal{I}^{\mathcal{F}}_{t,b}$	Mini-batch indices from ${\cal R}$ and ${\cal F}$ step $(t,b)$
B	Batch size	$\eta$	Learning rate
$\mathrm{sgn}(\cdot)$	Sign function returning $\pm 1$	$\alpha$	Weight in NegGrad balancing reta and forget contributions
$\widehat{m{\epsilon}}_{j,r}^{(t,b)}$	SAM perturbation applied to $\mathbf{w}_{j,r}^{(t,b)}$	au, ho	Perturbation radius in theory and practice used in SAM/ASAM
δ	Perturbation term: $\delta = \hat{\epsilon}_{j,r}^{(t,b)}$ for SAM and $0$ for SGD	$ abla_{oldsymbol{arphi}_i},  abla_{oldsymbol{\xi}_i}$	Gradient contributions for the sig and noise in NegGrad updates
$\Delta_{ ext{epoch}}^{ ext{SAM}} \kappa_{j,r}$	Per-epoch change of $\kappa_{j,r}$ under SAM	$\Delta_{ ext{epoch}}^{ ext{SGD}} \kappa_{j,r}$	Per-epoch change of $\kappa_{j,r}$ under SGD
$Acc(\theta, \mathcal{D})$	Classification accuracy model on dataset; $\theta$ , $\mathcal{D}$ are abbreviated terms in $\mathrm{Acc}()$	$ToW(f_{\mathcal{U}})$	"Tug-of-war" metric combining retain, forget and test accuracies
$\mathcal{D}$	data distribution	$\mathcal{F}_{high}$	Forget sets of high memorization difficulty; same for mid, low
$mem(\mathcal{A}, \mathcal{S}, i)$	Memorization score: $\Pr[f(\mathcal{S}) = y_i] - \Pr[f(\mathcal{S} \setminus i) = y_i]$	$\mathcal{S}\setminus i$	Training set $S$ with sample $i$ removed
$oldsymbol{\phi}_i$	Feature embedding of sample $i$ used in entanglement analysis	$oldsymbol{\mu}_{\mathcal{R}},oldsymbol{\mu}_{\mathcal{F}},oldsymbol{\mu}$	Mean embeddings of retain set, forget set and all data

Continued on next page

Table 5: (Continued)

Symbol	Meaning / Notes	Symbol	Meaning / Notes
$\mathrm{E}^{\mathrm{All}}_{\mathrm{Var}}(\mathcal{R},\mathcal{F},f)$	Variance-based entanglement measure between $\mathcal{R}$ and $\mathcal{F}$ , given model $f$	$\mathrm{E}^{\mathrm{Cls}}_{\mathrm{Var}}$	Class-wise version of the variance-based entanglement
$\mathcal{E}_{W_p}^{\mathrm{All}},\mathcal{E}_{W_p}^{\mathrm{Cls}}$	Geometry-aware entanglement measures based on Wasserstein distance (all/class-wise)	$\mathcal{A},\mathcal{U}$	Training algorithm (e.g. SGD, SAM) and unlearning algorithm (e.g. NegGrad, RL, with default SGD optimization, can be used with SAM)
$\kappa_{j,r}^{(0,0)}$	Initial signal coefficient for filter $(j,r)$	$\zeta_{j,r,i}^{(0,0)}$	Initial noise coefficient for sample $i$ on filter $\left(j,r\right)$
$ \mathcal{F} , \mathcal{R} $	Cardinalities of the forget and retain sets, which is size in our work	n	Total number of samples ( $ \mathcal{S} $ )
$\mathcal{D}_{ ext{test}}$	Test dataset used for evaluation	$\alpha^{\mathrm{SGD}},\alpha^{\mathrm{SAM}}$	$\alpha$ weight coeff for SGD and SAM, respectively

## DETAILED FORMULATIONS AND PROOFS

We prove our theorems and lemmas based on previous theoretical results in Kou et al. (2023); Chen et al. (2023). Specifically, we prove that with additional yet necessary conditions for effective unlearning, the final test errors can be preserved, while we identify and characterize the changed internal dynamics. We begin by expanding and restating  $\kappa, \zeta$  update rule for NegGrad in Eq. 7:

$$\kappa_{j,r}^{(t,b+1)} - \kappa_{j,r}^{(t,b)} = -\frac{\eta \|\boldsymbol{\varphi}\|_{2}^{2}}{Bm} \left[ \alpha \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{R}}} \ell_{i}^{\prime(t,b)} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t,b)} + \Delta, y_{i} \boldsymbol{\varphi} \rangle) - (1 - \alpha) \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{F}}} \ell_{i}^{\prime(t,b)} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t,b)} + \Delta, y_{i} \boldsymbol{\varphi} \rangle) \right],$$

$$\overline{\zeta}_{j,r}^{(t,b+1)} - \overline{\zeta}_{j,r}^{(t,b)} = -\frac{\eta (P - 1)^{2}}{Bm} \left[ \alpha \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{R}}} \|\boldsymbol{\xi}_{i}\|_{2}^{2} \ell_{i}^{\prime(t,b)} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t,b)} + \Delta, \boldsymbol{\xi}_{i} \rangle) \cdot \mathbb{1}(y_{i} = j) - (1 - \alpha) \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{F}}} \|\boldsymbol{\xi}_{i}\|_{2}^{2} \ell_{i}^{\prime(t,b)} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t,b)} + \Delta, \boldsymbol{\xi}_{i} \rangle) \cdot \mathbb{1}(y_{i} = j) \right],$$

$$\underline{\zeta}_{j,r}^{(t,b+1)} - \underline{\zeta}_{j,r}^{(t,b)} = +\frac{\eta (P - 1)^{2}}{Bm} \left[ \alpha \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{R}}} \|\boldsymbol{\xi}_{i}\|_{2}^{2} \ell_{i}^{\prime(t,b)} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t,b)} + \Delta, \boldsymbol{\xi}_{i} \rangle) \cdot \mathbb{1}(y_{i} \neq j) - (1 - \alpha) \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{F}}} \|\boldsymbol{\xi}_{i}\|_{2}^{2} \ell_{i}^{\prime(t,b)} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t,b)} + \Delta, \boldsymbol{\xi}_{i} \rangle) \cdot \mathbb{1}(y_{i} \neq j) \right],$$

$$(11)$$

where  $\Delta = \widehat{\epsilon}_{j,r}^{(t,b)}$  for SAM and 0 for SGD,  $\zeta_{j,r}^{(t,b)}$  is split into  $\overline{\zeta}_{j,r}^{(t,b)} := \zeta_{j,r}^{(t,b)} \mathbb{1}(\zeta_{j,r}^{(t,b)} \geq 0)$  and  $\underline{\zeta}_{j,r}^{(t,b)} := \zeta_{j,r}^{(t,b)} \mathbb{1}(\zeta_{j,r}^{(t,b)} \leq 0)$  based on label agreement. We summarize several reasonable assumptions from previous work in addition to our conditions which ensure unlearning to progress:

**Assumption D.1** Suppose there exists a sufficiently large constant C, such that the following hold:

1. Sufficiently large dimension  $d: d \ge C \max\{n\sigma_p^{-2} \|\varphi\|_2^2 \log(T^*), n^2 \log(nm/\delta)(\log(T^*))^2\}$ , for some  $T^* = \Omega(\eta^{-1}Bmd^{-1}P^{-2}\sigma_n^{-2}).$ 

- 2. The size of S and the CNN width satisfy  $n \ge C \log(m/\delta)$ ,  $m \ge C \log(n/\delta)$ .
- 3. The signal strength satisfies  $\|\varphi\|_2^2 \ge C\sigma_p^2 \log(n/\delta)$ .
- 4. For the Gaussian noise initialization,  $\sigma_0 \leq (C \max\{\sigma_v d/\sqrt{n}, \sqrt{\log(m/\delta)} \cdot \|\varphi\|_2\})^{-1}$ .
- 5. The learning rate  $\eta$  satisfies  $\eta \leq (C \max\{\sigma_v^2 d^{3/2}/(n^2 m \sqrt{\log(n/\delta)}), \sigma_v^2 d/n\})^{-1}$ .
- 6. Assume cross-entropy loss:  $\ell(z) = \log(1 + \exp(-z)) \Longrightarrow \ell' = -1/(1 + \exp(z))$ .
- 7. Assume ReLU activation.

8. Assume all clean labels and  $\mathcal{F}$  signals do not dominate:  $\alpha \geq |\mathcal{R}|/(|\mathcal{F}| + |\mathcal{R}|) := \beta > 0.5$ .

We then obtain several proven quantities from previous work, which are achieved during pretraining and can be leveraged at the start of unlearning:

- $\sum_{i=1}^n \overline{\zeta}_{j,r,i}^{(t)}/\kappa_{j',r'}^{(t)} = \Theta(SNR^{-2})$ , for the signal-to-noise ratio  $SNR = \frac{\|\varphi\|_2}{(P-1)\sigma_n\sqrt{d}}$ .
- $\sum_{i=1}^n \overline{\zeta}_{j,r,i}^{(t)} = \Omega(n) = O(n\log(T^*)) = \widetilde{\Theta}(n)$ , for some  $T^* = \Omega(\eta^{-1}Bmd^{-1}P^{-2}\sigma_p^{-2})$ .
- $\max_{j,r,i} |\underline{\zeta}_{j,r,i}^{(t)}| = \max\{O(\sqrt{\log(mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d}), O(\sqrt{\log(n/\delta)} \log(T^*) \cdot n/\sqrt{d})\}.$
- $\kappa_{i,r}^{(T^*)} = \Theta(\widehat{\kappa})$ , where  $\widehat{\kappa} = n \cdot \text{SNR}^2$ .

## D.1 PROOF TO THEOREM 3.2

Under NegGrad, we want to predict retain samples in  $\mathcal{R}$  correctly while we count correct predictions in  $\mathcal{F}$  as errors, yielding same bounds for  $\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}(yf(\mathbf{W}^{(t)},\mathbf{x})\leq 0)$  and  $\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{F}}(yf(\mathbf{W}^{(t)},\mathbf{x})>0)$  based on inverse objectives. However, when considering the test error on the model that is jointly updated by gradient descent on  $\mathcal{R}$  and gradient ascent on  $\mathcal{F}$ , we still measure the error rate by wrong predictions. In other words, fitting forget samples will reduce the generalization performance. We can decompose the test error as follows:

$$\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{D}}\left(y \neq \operatorname{sign}\left(f\left(\mathbf{W}^{(t)},\mathbf{x}\right)\right)\right) = \mathbb{P}_{(\mathbf{x},y)\sim\mathcal{D}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right) \leq 0\right) \\
= \mathbb{P}_{(\mathbf{x},y)\sim\mathcal{D}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right) \leq 0, (\mathbf{x},y) \in \mathcal{R}\right) + \mathbb{P}_{(\mathbf{x},y)\sim\mathcal{D}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right) \leq 0, (\mathbf{x},y) \in \mathcal{F}\right) \\
= \beta \cdot \mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right) \leq 0\right) + (1-\beta) \cdot \mathbb{P}_{(\mathbf{x},y)\sim\mathcal{F}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right) \leq 0\right) \\
= \beta \cdot \mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right) \leq 0\right) + (1-\beta) \cdot \left(1-\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{F}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right) > 0\right)\right).$$
(12)

Note that in practice,  $\mathcal{R}$  and  $\mathcal{F}$  come from training set  $\mathcal{S}$ . During inference and evaluation, we convert the data augmentations of  $\mathcal{R}$ ,  $\mathcal{F}$  to test transforms, thus measuring proxy-test errors on  $\mathcal{R}$ -like samples. To bound the test error, first decompose  $yf(\mathbf{W}^{(t)},\mathbf{x})$  into signal and noise learning of both positive and negative classes, considering  $\Delta=0$  for SGD:

$$yf\left(\mathbf{W}^{(t)}, \mathbf{x}\right) = \frac{1}{m} \sum_{j,r} yj \left[\sigma\left(\left\langle \mathbf{w}_{j,r}^{(t)}, y\varphi\right\rangle\right) + \sigma\left(\left\langle \mathbf{w}_{j,r}^{(t)}, \xi\right\rangle\right)\right]$$

$$= \frac{1}{m} \sum_{r} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, y\varphi\right\rangle\right) + (P-1)\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \xi\right\rangle\right)\right]$$

$$- \frac{1}{m} \sum_{r} \left[\sigma\left(\left\langle \mathbf{w}_{-y,r}^{(t)}, y\varphi\right\rangle\right) + (P-1)\sigma\left(\left\langle \mathbf{w}_{-y,r}^{(t)}, \xi\right\rangle\right)\right].$$
(13)

**Remark D.2** The following proof process for bounding  $\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}(yf(\mathbf{W}^{(t)},\mathbf{x}))$  comes from Kou et al. (2023). We include it here for readability, since we will leverage the results when combining  $\mathcal{R}$  and  $\mathcal{F}$  in the end, as well as make adaptations for proving Theorem 3.3. Our results benefit from previous work as we consider the unlearning process as an extension of the second stage in Chen et al. (2023).

We begin by two lemmas that bound the signal, noise norm, and the related inner products:

**Lemma D.3** (Lemma B.4 in Kou et al. (2023)). Suppose that  $\delta > 0$  and  $d = \Omega(\log(6n/\delta))$ . Then with probability at least  $1 - \delta$ ,

$$\sigma_p^2 d/2 \le \|\boldsymbol{\xi}_i\|_2^2 \le 3\sigma_p^2 d/2,$$
$$|\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_{i'} \rangle| \le 2\sigma_p^2 \cdot \sqrt{d \log(6n^2/\delta)},$$
$$|\langle \boldsymbol{\xi}_i, \boldsymbol{\varphi} \rangle| \le \|\boldsymbol{\varphi}\|_2 \sigma_p \cdot \sqrt{2 \log(6n/\delta)}.$$

for all  $i, i' \in [n]$ .

**Lemma D.4** (Lemma B.5 in Kou et al. (2023)). Suppose that  $d = \Omega(\log(mn/\delta)), m = \Omega(\log(1/\delta))$ . Then with probability at least  $1 - \delta$ ,

$$\sigma_0^2 d/2 \le \left\| \mathbf{w}_{j,r}^{(0,0)} \right\|_2^2 \le 3\sigma_0^2 d/2,$$

$$\left| \left\langle \mathbf{w}_{j,r}^{(0,0)}, \boldsymbol{\varphi} \right\rangle \right| \le \sqrt{2 \log(12m/\delta)} \cdot \sigma_0 \|\boldsymbol{\varphi}\|_2,$$

$$\left| \left\langle \mathbf{w}_{j,r}^{(0,0)}, \boldsymbol{\xi}_i \right\rangle \right| \le 2\sqrt{\log(12mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d},$$

for all  $r \in [m], j \in \{\pm 1\}$  and  $i \in [n]$ . Moreover,

$$\sigma_0 \|\varphi\|_2 / 2 \le \max_{r \in [m]} j \cdot \left\langle \mathbf{w}_{j,r}^{(0,0)}, \varphi \right\rangle \le \sqrt{2 \log(12m/\delta)} \cdot \sigma_0 \|\varphi\|_2,$$
  
$$\sigma_0 \sigma_p \sqrt{d} / 4 \le \max_{r \in [m]} j \cdot \left\langle \mathbf{w}_{j,r}^{(0,0)}, \boldsymbol{\xi}_i \right\rangle \le 2 \sqrt{\log(12mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d},$$

for all  $j \in \{\pm 1\}$  and  $i \in [n]$ .

Plug in the weight update decomposition in Eq. 2, we can first bound the inner product for j = y:

$$\begin{split} \left\langle \mathbf{w}_{y,r}^{(t)}, y\varphi \right\rangle &= \left\langle \mathbf{w}_{y,r}^{(0)}, y\varphi \right\rangle + \kappa_{y,r}^{(t)} \\ &+ \frac{1}{P-1} \sum_{i=1}^{n} \overline{\zeta}_{y,r,i}^{(t)} \left\| \boldsymbol{\xi}_{i} \right\|_{2}^{-2} \left\langle \boldsymbol{\xi}_{i}, y\varphi \right\rangle + \frac{1}{P-1} \sum_{i=1}^{n} \underline{\zeta}_{y,r,i}^{(t)} \left\| \boldsymbol{\xi}_{i} \right\|_{2}^{-2} \left\langle \boldsymbol{\xi}_{i}, y\varphi \right\rangle \\ &\geq -\sqrt{2 \log(12m/\delta)} \cdot \sigma_{0} \|\varphi\|_{2} + \kappa_{y,r}^{(t)} \\ &- \frac{\sqrt{2 \log(6n/\delta)}}{P-1} \cdot \sigma_{p} \|\varphi\|_{2} \cdot \left(\sigma_{p}^{2}d/2\right)^{-1} \left[ \sum_{i=1}^{n} \overline{\zeta}_{y,r,i}^{(t)} + \sum_{i=1}^{n} \left| \underline{\zeta}_{y,r,i}^{(t)} \right| \right] \\ &= -\Theta\left(\sqrt{\log(m/\delta)}\sigma_{0} \|\varphi\|_{2}\right) + \kappa_{y,r}^{(t)} - \Theta\left(\sqrt{\log(n/\delta)} \left(P\sigma_{p}d\right)^{-1} \|\varphi\|_{2}\right) \cdot \Theta\left(SNR^{-2}\right) \cdot \kappa_{y,r}^{(t)} \\ &= -\Theta\left(\sqrt{\log(m/\delta)} \left(\sigma_{p}d\right)^{-1} \sqrt{n} \|\varphi\|_{2}\right) + \left[1 - \Theta\left(\sqrt{\log(n/\delta)} \cdot P\sigma_{p}/\|\varphi\|_{2}\right) \right] \kappa_{y,r}^{(t)} \\ &= \Theta\left(\kappa_{y,r}^{(t)}\right), \end{split}$$

where the inequality is by Lemma D.3 and Lemma D.4; the second equality is obtained by plugging in the coefficient orders we summarized at the beginning of the section; the third equality is by  $\sigma_0 \leq C^{-1}(\sigma_p d)^{-1}\sqrt{n}$  in Assumption D.1 and SNR =  $\|\varphi\|_2/((P-1)\sigma_p\sqrt{d})$ . The fourth equality is by  $\kappa_{j,r}^{(t)} = \Theta(\widehat{\kappa})$ , where  $\widehat{\kappa} = n \cdot \text{SNR}^2$ . Also  $\sqrt{\log(n/\delta)} \cdot \sigma_p/\|\varphi\|_2 \leq 1/\sqrt{C}$  and  $\sqrt{\log(m/\delta)}(\sigma_p d)^{-1}\sqrt{n}\|\varphi\|_2/\widehat{\kappa} = \sqrt{\log(m/\delta)}\sigma_p/(\sqrt{n}\|\varphi\|_2) \leq \sqrt{\log(m/\delta)/n} \cdot 1/(\sqrt{C\log(n/\delta)}) \leq 1/(C\sqrt{\log(n/\delta)})$  holds by  $\|\varphi\|_2^2 \geq C \cdot \sigma_p^2 \log(n/\delta)$  and  $n \geq C\log(m/\delta)$  in Assumption D.1, so for sufficiently large constant C the equality holds. Similarly, we can show that  $\langle \mathbf{w}_{-y,r}^{(t)}, y\varphi \rangle = -\Theta(\kappa_{y,r}^{(t)}) < 0$  for  $j \neq y$ .

Next denote  $g(\boldsymbol{\xi})$  as  $\sum_{r} \sigma(\langle \mathbf{w}_{-y,r}^{(t)}, \boldsymbol{\xi} \rangle)$ . Since  $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$ , we can leverage the Gaussian concentration bound for  $x \geq 0$ :

$$\mathbb{P}(g(\boldsymbol{\xi}) - \mathbb{E}g(\boldsymbol{\xi}) \ge x) \le \exp\left(-\frac{cx^2}{\sigma_p^2 \|g\|_{\text{Lip}}^2}\right),\tag{15}$$

where c is a constant. To calculate the Lipschitz norm, we have

$$|g(\boldsymbol{\xi}) - g(\boldsymbol{\xi}')| = \left| \sum_{r=1}^{m} \sigma\left( \left\langle \mathbf{w}_{-y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) - \sum_{r=1}^{m} \sigma\left( \left\langle \mathbf{w}_{-y,r}^{(t)}, \boldsymbol{\xi}' \right\rangle \right) \right|$$

$$\leq \sum_{r=1}^{m} \left| \sigma\left( \left\langle \mathbf{w}_{-y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left( \left\langle \mathbf{w}_{-y,r}^{(t)}, \boldsymbol{\xi}' \right\rangle \right) \right|$$

$$\leq \sum_{r=1}^{m} \left| \left\langle \mathbf{w}_{-y,r}^{(t)}, \boldsymbol{\xi} - \boldsymbol{\xi}' \right\rangle \right| \leq \sum_{r=1}^{m} \left\| \mathbf{w}_{-y,r}^{(t)} \right\|_{2} \cdot \|\boldsymbol{\xi} - \boldsymbol{\xi}'\|_{2}.$$
(16)

The first inequality is by triangle inequality; the second inequality is by the property of ReLU; the last inequality is by Cauchy-Schwartz inequality. Therefore, we have  $\|g\|_{\mathrm{Lip}} \leq \sum_{r=1}^m \|\mathbf{w}_{-y,r}^{(t)}\|_2$ , and since  $\langle \mathbf{w}_{-y,r}^{(t)}, \boldsymbol{\xi} \rangle \sim \mathcal{N}(0, \|\mathbf{w}_{-y,r}^{(t)}\|_2^2 \sigma_p^2)$ , we can get

$$\mathbb{E}g(\boldsymbol{\xi}) = \sum_{r=1}^{m} \mathbb{E}\sigma\left(\left\langle \mathbf{w}_{-y,r}^{(t)}, \boldsymbol{\xi} \right\rangle\right) = \sum_{r=1}^{m} \frac{\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2} \sigma_{p}}{\sqrt{2\pi}} = \frac{\sigma_{p}}{\sqrt{2\pi}} \sum_{r=1}^{m} \left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}.$$
 (17)

Then, we seek to upper bound the 2-norm of  $\mathbf{w}_{i\,r}^{(t)}$ . First we have

$$\left\| \sum_{i=1}^{n} \zeta_{j,r,i}^{(t)} \cdot \|\boldsymbol{\xi}_{i}\|_{2}^{-2} \cdot \boldsymbol{\xi}_{i} \right\|_{2}^{2} \\
= \sum_{i=1}^{n} \zeta_{j,r,i}^{(t)} \cdot \|\boldsymbol{\xi}_{i}\|_{2}^{-2} + 2 \underbrace{\sum_{1 \leq i_{1} < i_{2} \leq n} \zeta_{j,r,i_{1}}^{(t)} \zeta_{j,r,i_{2}}^{(t)} \cdot \|\boldsymbol{\xi}_{i_{1}}\|_{2}^{-2} \|\boldsymbol{\xi}_{i_{2}}\|_{2}^{-2} \cdot \langle \boldsymbol{\xi}_{i_{1}}, \boldsymbol{\xi}_{i_{2}} \rangle}_{\text{off-diagonal}} \\
\leq 4\sigma_{p}^{-2} d^{-1} \sum_{i=1}^{n} \zeta_{j,r,i}^{(t)}^{2} + 2 \underbrace{\sum_{1 \leq i_{1} < i_{2} \leq n} \left| \zeta_{j,r,i_{1}}^{(t)} \zeta_{j,r,i_{2}}^{(t)} \right| \cdot \left(16\sigma_{p}^{-4} d^{-2}\right) \cdot \left(2\sigma_{p}^{2} \sqrt{d \log (6n^{2}/\delta)}\right)}_{1 \leq 4\sigma_{p}^{-2} d^{-1} \sum_{i=1}^{n} \zeta_{j,r,i}^{(t)}^{2} + 32\sigma_{p}^{-2} d^{-3/2} \sqrt{\log (6n^{2}/\delta)} \left[ \left(\sum_{i=1}^{n} \left| \zeta_{j,r,i}^{(t)} \right| \right)^{2} - \sum_{i=1}^{n} \zeta_{j,r,i}^{(t)}^{2} \right] \\
= \Theta\left(\sigma_{p}^{-2} d^{-1}\right) \sum_{i=1}^{n} \zeta_{j,r,i}^{(t)}^{2} + \widetilde{\Theta}\left(\sigma_{p}^{-2} d^{-3/2}\right) \left(\sum_{i=1}^{n} \left| \zeta_{j,r,i}^{(t)} \right| \right)^{2} \\
\leq \left[ \Theta\left(\sigma_{p}^{-2} d^{-1} n^{-1}\right) + \widetilde{\Theta}\left(\sigma_{p}^{-2} d^{-3/2}\right) \right] \left(\sum_{i=1}^{n} \left| \zeta_{j,r,i}^{(t)} \right| + \sum_{i=1}^{n} \left| \zeta_{j,r,i}^{(t)} \right| \right)^{2} \\
\leq \Theta\left(\sigma_{p}^{-2} d^{-1} n^{-1}\right) \left(\sum_{i=1}^{n} \zeta_{j,r,i}^{(t)}\right)^{2}.$$
(18)

The first inequality is by Lemma D.3; for the second inequality we used the definition of  $\overline{\zeta}$ ,  $\underline{\zeta}$ ; for the second to last equation we plugged in coefficient orders. We can thus upper bound the 2-norm of  $\mathbf{w}_{i,r}^{(t)}$  as:

$$\left\| \mathbf{w}_{j,r}^{(t)} \right\|_{2} \leq \left\| \mathbf{w}_{j,r}^{(0)} \right\|_{2} + \kappa_{j,r}^{(t)} \cdot \left\| \varphi \right\|_{2}^{-1} + \frac{1}{P-1} \left\| \sum_{i=1}^{n} \zeta_{j,r,i}^{(t)} \cdot \left\| \boldsymbol{\xi}_{i} \right\|_{2}^{-2} \cdot \boldsymbol{\xi}_{i} \right\|_{2}$$

$$\leq \left\| \mathbf{w}_{j,r}^{(0)} \right\|_{2} + \kappa_{j,r}^{(t)} \cdot \left\| \varphi \right\|_{2}^{-1} + \Theta \left( P^{-1} \sigma_{p}^{-1} d^{-1/2} n^{-1/2} \right) \cdot \sum_{i=1}^{n} \overline{\zeta}_{j,r,i}^{(t)}$$

$$= \Theta \left( P^{-1} \sigma_{p}^{-1} d^{-1/2} n^{-1/2} \right) \cdot \sum_{i=1}^{n} \overline{\zeta}_{j,r,i}^{(t)},$$

$$(19)$$

 where the first inequality is due to the triangle inequality, and the equality is due to the following:

$$\frac{\kappa_{j,r}^{(t)} \cdot \|\varphi\|_{2}^{-1}}{\Theta\left(P^{-1}\sigma_{p}^{-1}d^{-1/2}n^{-1/2}\right) \cdot \sum_{i=1}^{n} \overline{\zeta}_{j,r,i}^{(t)}} = \Theta\left(P^{-1}\sigma_{p}d^{1/2}n^{1/2}\|\varphi\|_{2}^{-1}SNR^{2}\right) 
= \Theta\left(P^{-1}\sigma_{p}^{-1}d^{-1/2}n^{1/2}\|\varphi\|_{2}\right) = O(1),$$
(20)

based on the coefficient order  $\sum_{i=1}^n \overline{\zeta}_{j,r,i}^{(t)}/\kappa_{j,r}^{(t)} = \Theta(\mathrm{SNR}^{-2})$ , the definition of SNR, and the condition for d in Assumption D.1. Similarly,

$$\frac{\left\|\mathbf{w}_{j,r}^{(0)}\right\|_{2}}{\Theta\left(P^{-1}\sigma_{p}^{-1}d^{-1/2}n^{-1/2}\right)\cdot\sum_{i=1}^{n}\overline{\zeta}_{j,r,i}^{(t)}} = \frac{\Theta\left(\sigma_{0}\sqrt{d}\right)}{\Theta\left(P^{-1}\sigma_{p}^{-1}d^{-1/2}n^{-1/2}\right)\cdot\sum_{i=1}^{n}\overline{\zeta}_{j,r,i}^{(t)}} = O\left(P\sigma_{0}\sigma_{p}dn^{-1/2}\right) = O(1),$$
(21)

based on Lemma D.4, the coefficient order  $\sum_{i=1}^{n} \overline{\zeta}_{j,r,i}^{(t)} = \Omega(n)$ , and the condition for  $\sigma_0$  in Assumption D.1. Then we can give an analysis of the following key component:

$$\frac{\sum_{r} \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, y \boldsymbol{\varphi} \right\rangle\right)}{(P-1)\sigma_{p} \sum_{r=1}^{m} \left\| \mathbf{w}_{-y,r}^{(t)} \right\|_{2}} \ge \frac{\Theta\left(\sum_{r} \kappa_{y,r}^{(t)}\right)}{\Theta\left(d^{-1/2} n^{-1/2}\right) \cdot \sum_{r,i} \overline{\zeta}_{-y,r,i}^{(t)}} \\
= \Theta\left(d^{1/2} n^{1/2} \text{SNR}^{2}\right) = \Theta\left(n^{1/2} \|\boldsymbol{\varphi}\|_{2}^{2} / (P^{2} \sigma_{p}^{2} d^{1/2})\right). \tag{22}$$

Then for  $\|\varphi\|_2 \ge C_1^{1/4} n^{-1/4} P \sigma_p d^{1/4}$  for some large constant  $C_1$ , we have

$$\sum_{r} \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, y\varphi \right\rangle\right) - \frac{(P-1)\sigma_{p}}{\sqrt{2\pi}} \sum_{r=1}^{m} \left\| \mathbf{w}_{-y,r}^{(t)} \right\|_{2} > 0.$$
 (23)

Upper bound. Now plug in previous results to obtain

$$\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right)\leq0\right)\leq\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left((P-1)\sum_{r}\sigma\left(\left\langle\mathbf{w}_{-y,r}^{(t)},\boldsymbol{\xi}\right\rangle\right)\geq\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\boldsymbol{\varphi}\right\rangle\right)\right)$$

$$=\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(g(\boldsymbol{\xi})-\mathbb{E}g(\boldsymbol{\xi})\geq1/(P-1)\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\boldsymbol{\varphi}\right\rangle\right)-\frac{\sigma_{p}}{\sqrt{2\pi}}\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}\right)$$

$$\leq\exp\left[-\frac{c\left(1/(P-1)\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\boldsymbol{\varphi}\right\rangle\right)-\left(\sigma_{p}/\sqrt{2\pi}\right)\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}\right)^{2}}{\sigma_{p}^{2}\left(\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}\right)^{2}}\right]$$

$$=\exp\left[-c\left(\frac{\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\boldsymbol{\varphi}\right\rangle\right)}{(P-1)\sigma_{p}\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}}-1/\sqrt{2\pi}\right)^{2}\right]$$

$$\leq\exp(c/2\pi)\exp\left(-0.5c\left(\frac{\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\boldsymbol{\varphi}\right\rangle\right)}{(P-1)\sigma_{p}\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}}\right)^{2}\right).$$
(24)

The second inequality is by Eq. 23 and plugging  $\|g\|_{\text{Lip}} \le \sum_{r=1}^m \|\mathbf{w}_{-y,r}^{(t)}\|_2$  into Eq. 15; the third inequality is due to  $(s-t)^2 \ge s^2/2 - t^2, \forall s,t \ge 0$ . And from Eq. 22 and Eq. 24 we have

$$\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right)\leq 0\right) \leq \exp(c/2\pi) \exp\left(-0.5c\left(\frac{\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\varphi\right\rangle\right)}{(P-1)\sigma_{p}\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}}\right)^{2}\right)$$

$$= \exp\left(\frac{c}{2\pi} - \frac{n\|\varphi\|_{2}^{4}}{C(P-1)^{4}\sigma_{p}^{4}d}\right)$$

$$\leq \exp\left(-\frac{n\|\varphi\|_{2}^{4}}{2C_{1}(P-1)^{4}\sigma_{p}^{4}d}\right)$$

$$= \exp\left(-\frac{n\|\varphi\|_{2}^{4}}{C_{2}(P-1)^{4}\sigma_{p}^{4}d}\right) = \epsilon,$$
(25)

where C = O(1); the last inequality holds if we choose  $C_1 \ge cC/\pi$ ; the last equality holds if we choose  $C_2$  as 2C.

For the forget set  $\mathcal{F}$ , we thus have

$$\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{F}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right)>0\right)\leq\epsilon.$$
(26)

**Lower bound.** Without loss of generality, let  $\sum_{r} \kappa_{1,r}^{(t)} = \max \left\{ \sum_{r} \kappa_{1,r}^{(t)}, \sum_{r} \kappa_{-1,r}^{(t)} \right\}$ . Denote  $\mathbf{v} = \lambda \cdot \sum_{i} \mathbbm{1}(y_i = 1) \boldsymbol{\xi}_i$ , where  $\lambda = C_7 \mathrm{SNR}^2 = C_7 \|\boldsymbol{\varphi}\|_2^2 / \left( (P-1)^2 \sigma_p^2 d \right)$  and  $C_7$  is a sufficiently large constant. Since ReLU is convex, we have

$$\sigma\left(\left\langle \mathbf{w}_{1,r}^{(t)}, \boldsymbol{\xi} + \mathbf{v} \right\rangle\right) - \sigma\left(\left\langle \mathbf{w}_{1,r}^{(t)}, \boldsymbol{\xi} \right\rangle\right) \ge \sigma'\left(\left\langle \mathbf{w}_{1,r}^{(t)}, \boldsymbol{\xi} \right\rangle\right) \left\langle \mathbf{w}_{1,r}^{(t)}, \mathbf{v} \right\rangle,$$

$$\sigma\left(\left\langle \mathbf{w}_{1,r}^{(t)}, -\boldsymbol{\xi} + \mathbf{v} \right\rangle\right) - \sigma\left(\left\langle \mathbf{w}_{1,r}^{(t)}, -\boldsymbol{\xi} \right\rangle\right) \ge \sigma'\left(\left\langle \mathbf{w}_{1,r}^{(t)}, -\boldsymbol{\xi} \right\rangle\right) \left\langle \mathbf{w}_{1,r}^{(t)}, \mathbf{v} \right\rangle.$$
(27)

Summing the above two, we have that almost surely for all  $\xi$ 

$$\sigma\left(\left\langle \mathbf{w}_{1,r}^{(t)}, \boldsymbol{\xi} + \mathbf{v} \right\rangle\right) - \sigma\left(\left\langle \mathbf{w}_{1,r}^{(t)}, \boldsymbol{\xi} \right\rangle\right) + \sigma\left(\left\langle \mathbf{w}_{1,r}^{(t)}, -\boldsymbol{\xi} + \mathbf{v} \right\rangle\right) - \sigma\left(\left\langle \mathbf{w}_{1,r}^{(t)}, -\boldsymbol{\xi} \right\rangle\right) \\
\geq \left\langle \mathbf{w}_{1,r}^{(t)}, \mathbf{v} \right\rangle \\
\geq \lambda \left[ \sum_{y_{i}=1} \overline{\zeta}_{1,r,i}^{(t)} - 2n\sqrt{\log(12mn/\delta)} \cdot \sigma_{0}\sigma_{p}\sqrt{d} - 5n^{2}\alpha\sqrt{\log(6n^{2}/\delta)/d} \right], \tag{28}$$

where the last inequality is by Lemma C.3 in Kou et al. (2023) and Lemma D.4. Additionally, since ReLU is a Liptchitz, we also have that

$$\sigma\left(\left\langle \mathbf{w}_{-1,r}^{(t)}, \boldsymbol{\xi} + \mathbf{v} \right\rangle\right) - \sigma\left(\left\langle \mathbf{w}_{-1,r}^{(t)}, \boldsymbol{\xi} \right\rangle\right) + \sigma\left(\left\langle \mathbf{w}_{-1,r}^{(t)}, -\boldsymbol{\xi} + \mathbf{v} \right\rangle\right) - \sigma\left(\left\langle \mathbf{w}_{-1,r}^{(t)}, -\boldsymbol{\xi} \right\rangle\right) \\
\leq 2 \left|\left\langle \mathbf{w}_{-1,r}^{(t)}, \mathbf{v} \right\rangle\right| \\
\leq 2\lambda \left[\sum_{y_i=1} \underline{\zeta}_{-1,r,i}^{(t)} + 2n\sqrt{\log(12mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d} + 5n^2 \alpha \sqrt{\log(6n^2/\delta)/d}\right].$$
(29)

Therefore, by plugging Eq. 28 and Eq. 29, we have that

1190 
$$g(\boldsymbol{\xi} + \mathbf{v}) - g(\boldsymbol{\xi}) + g(-\boldsymbol{\xi} + \mathbf{v}) - g(-\boldsymbol{\xi})$$
1191  $\geq \lambda \left[ \sum_{r} \sum_{y_i=1} \overline{\zeta}_{1,r,i}^{(t)} - 6nm\sqrt{\log(12mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d} - 15mn^2 \alpha \sqrt{\log(6n^2/\delta)/d} \right]$ 
1194  $\geq (\lambda/2) \cdot \sum_{r} \sum_{y_i=1} \overline{\zeta}_{1,r,i}^{(t)}$ 
1196  $\geq \lambda/2 \cdot \Theta\left(\mathrm{SNR}^{-2}\right) \sum_{r} \kappa_{1,r}^{(t)}$ 
1198  $\geq 4C_6 \sum_{r} \kappa_{1,r}^{(t)}$ ,
1200

where the second inequality is by Lemma D.1 in Kou et al. (2023) and Assumption D.1; the third inequality is by  $\sum_{i=1}^{n} \overline{\zeta}_{j,r,i}^{(t)}/\kappa_{j',r'}^{(t)} = \Theta(\mathrm{SNR}^{-2})$ . Finally, it is worth noting that the norm

$$\|\mathbf{v}\|_{2} = \left\|\lambda \cdot \sum_{i} \mathbb{1}\left(y_{i} = 1\right) \boldsymbol{\xi}_{i}\right\|_{2} = \Theta\left(\sqrt{\frac{n\|\boldsymbol{\varphi}\|_{2}^{4}}{P^{4}\sigma_{p}^{4}d}}\right) \leq 0.06\sigma_{p}.$$
(31)

where the last inequality is by condition  $\|\varphi\|_2 \le C_3 d^{1/4} n^{-1/4} P \sigma_p$  with sufficiently large  $C_3$ . Then we present a Lemma which bounds the Total Variation (TV) distance between two Gaussian with the same covariance matrix.

**Lemma D.5** (Proposition 2.1 by Devroye et al. (2018)). The TV distance between  $\mathcal{N}\left(0, \sigma_p^2 \mathbf{I}_d\right)$  and  $\mathcal{N}\left(\mathbf{v}, \sigma_p^2 \mathbf{I}_d\right)$  is smaller than  $\|\mathbf{v}\|_2 / 2\sigma_p$ .

Finally, we can prove the lower bound for  $\mathcal{R}$ :

$$\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right)\leq0\right) \\
=\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(\sum_{r}\sigma\left(\left\langle\mathbf{w}_{-y,r}^{(t)},\boldsymbol{\xi}\right\rangle\right)-\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},\boldsymbol{\xi}\right\rangle\right)\geq\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\varphi\right\rangle\right)-\sum_{r}\sigma\left(\left\langle\mathbf{w}_{-y,r}^{(t)},y\varphi\right\rangle\right)\right) \\
\geq0.5\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(\left|\sum_{r}\sigma\left(\left\langle\mathbf{w}_{-y,r}^{(t)},\boldsymbol{\xi}\right\rangle\right)-\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},\boldsymbol{\xi}\right\rangle\right)\right|\geq C_{6}\max\left\{\sum_{r}\kappa_{1,r}^{(t)},\sum_{r}\kappa_{-1,r}^{(t)}\right\}\right), \tag{32}$$

where  $C_6$  is a constant, the inequality holds since if  $|\sum_r \sigma(\langle \mathbf{w}_{1,r}^{(t)}, \boldsymbol{\xi} \rangle) - \sum_r \sigma(\langle \mathbf{w}_{-1,r}^{(t)}, \boldsymbol{\xi} \rangle)|$  is too large, we can always pick a corresponding y given  $\boldsymbol{\xi}$  to make a wrong prediction.

Let  $g(\boldsymbol{\xi}) = \sum_r \sigma(\langle \mathbf{w}_{1,r}^{(t)}, \boldsymbol{\xi} \rangle) - \sum_r \sigma(\langle \mathbf{w}_{-1,r}^{(t)}, \boldsymbol{\xi} \rangle)$ , and denote the set  $\Omega := \{ \boldsymbol{\xi} \mid |g(\boldsymbol{\xi})| \geq C_6 \max\{\sum_r \kappa_{1.r}^{(t)}, \sum_r \kappa_{-1.r}^{(t)} \} \}$ . Thus we have

$$\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right)\leq 0\right)\geq 0.5\mathbb{P}(\Omega). \tag{33}$$

By Lemma 5.8 of Kou et al. (2023), we have that  $\sum_j [g(j\boldsymbol{\xi} + \mathbf{v}) - g(j\boldsymbol{\xi})] \ge 4C_6 \max_j \left\{ \sum_r \kappa_{j,r}^{(t)} \right\}$ . Therefore, by pigeonhole principle, one of  $[\boldsymbol{\xi}, -\boldsymbol{\xi}, \boldsymbol{\xi} + \mathbf{v}, -\boldsymbol{\xi} + \mathbf{v}]$  must belong to  $\Omega$ , thus  $\Omega \cup -\Omega \cup \Omega - \{\mathbf{v}\} \cup -\Omega - \{\mathbf{v}\} = \mathbb{R}^d$ . Therefore, at least one of  $\mathbb{P}(\Omega), \mathbb{P}(-\Omega), \mathbb{P}(\Omega - \{\mathbf{v}\}), \mathbb{P}(-\Omega - \{\mathbf{v}\})$  is greater than  $\frac{1}{4}$ . Note that  $\mathbb{P}(-\Omega) = \mathbb{P}(\Omega)$  and

$$|\mathbb{P}(\Omega) - \mathbb{P}(\Omega - \mathbf{v})| = \left| \mathbb{P}_{\boldsymbol{\xi} \sim \mathcal{N}\left(0, \sigma_{p}^{2} \mathbf{I}_{d}\right)}(\boldsymbol{\xi} \in \Omega) - \mathbb{P}_{\boldsymbol{\xi} \sim \mathcal{N}\left(\mathbf{v}, \sigma_{p}^{2} \mathbf{I}_{d}\right)}(\boldsymbol{\xi} \in \Omega) \right|$$

$$\leq \text{TV}\left(\mathcal{N}\left(0, \sigma_{p}^{2} \mathbf{I}_{d}\right), \mathcal{N}\left(\mathbf{v}, \sigma_{p}^{2} \mathbf{I}_{d}\right)\right)$$

$$\leq \frac{\|\mathbf{v}\|_{2}}{2\sigma_{p}} \leq 0.03,$$
(34)

where the first inequality is by the definition of TV distance, the second inequality is by Lemma D.5. Hence, we have that  $\mathbb{P}(\Omega) \geq \frac{1}{4} - 0.03 = 0.22$ , and plugging this into Eq. 33, we get

$$\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right)\leq 0\right)\geq 0.5\mathbb{P}(\Omega)=0.11\geq 0.1. \tag{35}$$

Like the upper bound, the derived lower bounds also applies to  $\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{F}}(yf(\mathbf{W}^{(t)},\mathbf{x})>0)$ . Hence, if  $\|\varphi\|_2 \geq C_1 d^{1/4} n^{-1/4} P \sigma_p$ ,

$$\mathcal{L}^{\text{test}}(\mathbf{W}^{T_{2}}, \mathcal{D}) = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left( y \neq \text{sign} \left( f\left(\mathbf{W}^{T_{2}}, \mathbf{x}\right) \right) \right)$$

$$= \beta \cdot \underbrace{\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{R}} \left( y f\left(\mathbf{W}^{T_{2}}, \mathbf{x}\right) \leq 0 \right)}_{\leq \epsilon_{\mathcal{R}}} + (1 - \beta) \cdot \left( 1 - \underbrace{\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{F}} \left( y f\left(\mathbf{W}^{T_{2}}, \mathbf{x}\right) > 0 \right)}_{\leq \epsilon_{\mathcal{F}}} \right)$$

$$\implies \lim_{\beta \to 1} \mathcal{L}^{\text{test}}(\mathbf{W}^{T_{2}}, \mathcal{D}) \leq \epsilon_{\mathcal{R}} = \epsilon.$$
(36)

On the other hand, when  $\beta \to 0.5$ , we have  $\lim_{\beta \to 0.5} \mathcal{L}^{\text{test}}(\mathbf{W}^{T_2}, \mathcal{D}) \le 0.5 + 0.5\epsilon_{\mathcal{R}} - 0.5\epsilon_{\mathcal{F}} = \epsilon$ . Depending on the size ratio of  $\mathcal{R}$  and  $\mathcal{F}$ ,  $\epsilon$  ranges from a very small constant to a minimally PAC-learnable threshold.

For harmful overfitting where  $\|\varphi\|_2 \leq C_3 d^{1/4} n^{-1/4} P \sigma_p$ ,

$$\mathcal{L}^{\text{test}}(\mathbf{W}^{T_{2}}, \mathcal{D}) = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left( y \neq \text{sign} \left( f\left(\mathbf{W}^{T_{2}}, \mathbf{x}\right) \right) \right)$$

$$= \beta \cdot \underbrace{\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{R}} \left( y f\left(\mathbf{W}^{T_{2}}, \mathbf{x}\right) \leq 0 \right)}_{\geq 0.1} + (1 - \beta) \cdot \left( 1 - \underbrace{\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{F}} \left( y f\left(\mathbf{W}^{T_{2}}, \mathbf{x}\right) > 0 \right)}_{\geq 0.1} \right)$$

$$\implies \lim_{\beta \to 1} \mathcal{L}^{\text{test}}(\mathbf{W}^{T_{2}}, \mathcal{D}) \geq 0.1.$$
(37)

On the other hand, when  $\beta \to 0.5$ , we have  $\lim_{\beta \to 0.5} \mathcal{L}^{\text{test}}(\mathbf{W}^{T_2}, \mathcal{D}) \ge 0.05$ .

## D.2 PROOF TO THEOREM 3.3

First we have the same decomposition for NegGrad:

$$\mathcal{L}^{\text{test}}(\mathbf{W}^{T_{2}}, \mathcal{D}) = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left( y \neq \text{sign} \left( f\left(\mathbf{W}^{(t)}, \mathbf{x}\right) \right) \right)$$

$$= \beta \cdot \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{R}} \left( y f\left(\mathbf{W}^{(t)}, \mathbf{x}\right) \leq 0 \right) + (1 - \beta) \cdot \left( 1 - \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{F}} \left( y f\left(\mathbf{W}^{(t)}, \mathbf{x}\right) > 0 \right) \right);$$

$$y f\left(\mathbf{W}^{(t)}, \mathbf{x}\right) = \frac{1}{m} \sum_{j,r} y j \left[ \sigma\left( \left\langle \mathbf{w}_{j,r}^{(t)}, y \varphi \right\rangle \right) + \sigma\left( \left\langle \mathbf{w}_{j,r}^{(t)}, \xi \right\rangle \right) \right]$$

$$= \frac{1}{m} \sum_{r} \left[ \sigma\left( \left\langle \mathbf{w}_{y,r}^{(t)}, y \varphi \right\rangle \right) + (P - 1) \sigma\left( \left\langle \mathbf{w}_{y,r}^{(t)}, \xi \right\rangle \right) \right]$$

$$- \frac{1}{m} \sum_{r} \left[ \sigma\left( \left\langle \mathbf{w}_{-y,r}^{(t)}, y \varphi \right\rangle \right) + (P - 1) \sigma\left( \left\langle \mathbf{w}_{-y,r}^{(t)}, \xi \right\rangle \right) \right].$$
(38)

However, note that for  $(\mathbf{x}, y) \sim \mathcal{F}$ , SAM gives up its denoising property. We first show this by proving Lemma 3.1.

## D.2.1 PROOF TO LEMMA 3.1

*Proof.* Consider extending Lemma D.5 in Chen et al. (2023) to the NegGrad setting by rewriting  $\left\langle \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_{k} \right\rangle$ . First we have the Frobenius norm upper bounded by the same quantity:

$$\|\nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{I}_{t,b}}(\mathbf{W}^{(t,b)})\|_{F} = \|\alpha \nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{I}_{t,b}^{\mathcal{R}}}(\mathbf{W}^{(t,b)}) - (1 - \alpha) \nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{I}_{t,b}^{\mathcal{F}}}(\mathbf{W}^{(t,b)})\|_{F}$$

$$\leq \alpha \|\nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{I}_{t,b}^{\mathcal{R}}}(\mathbf{W}^{(t,b)})\|_{F} + (1 - \alpha) \|\nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{I}_{t,b}^{\mathcal{F}}}(\mathbf{W}^{(t,b)})\|_{F}$$

$$= \|\nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{I}_{t,b}}(\mathbf{W}^{(t,b)})\|_{F} \leq 2\sqrt{2} P \sigma_{p} \sqrt{d/Bm},$$
(39)

where the first inequality comes from triangle inequality; the second equality holds because  $\mathcal{R}, \mathcal{F}$  are split from  $\mathcal{S}$  and come from the same  $\mathcal{D}$ , thus having the same gradient norm; the second inequality comes from the original bounds in Chen et al. (2023). Next we expand  $\left\langle \widehat{\epsilon}_{j,r}^{(t,b)}, \xi_k \right\rangle$  under NegGrad:

$$\left\langle \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_{k} \right\rangle = \frac{\tau}{mB} \left\| \nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{I}_{t,b}}(\mathbf{W}^{(t,b)}) \right\|_{F}^{-1} \sum_{i \in \mathcal{I}_{t,b}} \sum_{p \in [P]} \ell_{i}^{\prime(t)} j \cdot y_{i} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_{i,p} \rangle) \langle \mathbf{x}_{i,p}, \boldsymbol{\xi}_{k} \rangle$$

$$= \frac{\tau}{mB} \left\| \nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{I}_{t,b}}(\mathbf{W}^{(t,b)}) \right\|_{F}^{-1} \left[ \alpha \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{R}}} \sum_{p \in [P]} \ell_{i}^{\prime(t)} j \cdot y_{i} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_{i,p} \rangle) \langle \mathbf{x}_{i,p}, \boldsymbol{\xi}_{k} \rangle \right.$$

$$\left. - (1 - \alpha) \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{F}}} \sum_{p \in [P]} \ell_{i}^{\prime(t)} j \cdot y_{i} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_{i,p} \rangle) \langle \mathbf{x}_{i,p}, \boldsymbol{\xi}_{k} \rangle \right]. \tag{40}$$

Note that  $\langle \mathbf{x}_{i,p}, \boldsymbol{\xi}_k \rangle$  can be divided into three different terms:

$$|\langle \mathbf{x}_{i,p}, \boldsymbol{\xi}_{k} \rangle| = \begin{cases} \|\boldsymbol{\xi}_{k}\|_{2}^{2} \leq 3\sigma_{p}^{2}d/2, & \text{if } i = k, x_{k,p} = \boldsymbol{\xi}_{k} \\ |\langle \boldsymbol{\xi}_{i}, \boldsymbol{\xi}_{k} \rangle| \leq 2\sigma_{p}^{2}\sqrt{d\log(6n^{2}/\delta)}, & \text{if } i \neq k, x_{i,p} = \boldsymbol{\xi}_{i} \\ |\langle y_{i}\boldsymbol{\varphi}, \boldsymbol{\xi}_{k} \rangle| \leq \|\boldsymbol{\varphi}\|_{2} \sigma_{p}\sqrt{2\log(6n^{2}/\delta)}, & \text{if } x_{i,p} = y_{i}\boldsymbol{\varphi} \end{cases}$$
(41)

The upper bounds come from Lemma D.3. Based on Assumption D.1 and Lemma D.4 of Chen et al. (2023), the i = k term will dominate the upper bound and we can write

$$\left\langle \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_{k} \right\rangle \leq \frac{\tau}{mB \cdot 2\sqrt{2}P\sigma_{p}\sqrt{d/Bm}} \left[ -0.15\alpha(P-1)C_{1}\sigma_{p}^{2}d\mathbb{1}[k \in \mathcal{I}_{t,b}^{\mathcal{R}}] +0.15(1-\alpha)(P-1)C_{1}\sigma_{p}^{2}d\mathbb{1}[k \in \mathcal{I}_{t,b}^{\mathcal{F}}] \right]$$

$$(42)$$

Thus, when  $k \in \mathcal{I}_{t,b}^{\mathcal{R}}$ , we can preserve the original bound with additional  $\alpha$ :

$$\left\langle \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_{k} \right\rangle < -C \frac{\alpha \tau \sigma_{p} \sqrt{d}}{m \sqrt{B}}.$$
 (43)

Choosing  $\tau = \frac{m\sqrt{B}}{C_3\alpha P\sigma_p\sqrt{d}}$  will cancel with  $\left\langle \mathbf{w}_{j,r}^{(t)}, \boldsymbol{\xi}_k \right\rangle$  to deactivate the neuron. When  $k \in \mathcal{I}_{t,b}^{\mathcal{F}}$ , the entire  $\left\langle \mathbf{w}_{j,r}^{(t,b)} + \hat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \right\rangle$  will remain activated:

$$0 \le \left\langle \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \right\rangle < C \frac{(1-\alpha)\tau \sigma_p \sqrt{d}}{m\sqrt{B}} \Longrightarrow \left\langle \mathbf{w}_{j,r}^{(t,b)} + \widehat{\boldsymbol{\epsilon}}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \right\rangle \ge \left\langle \mathbf{w}_{j,r}^{(t,b)}, \boldsymbol{\xi}_k \right\rangle \ge 0. \tag{44}$$

This fundamentally differs SAM's behaviors towards unlearning  $\mathcal{F}$  from behaviors towards learning  $\mathcal{R}$  as how SGD differs from SAM. For gradient ascent on  $\mathcal{F}$  under NegGrad, we now know SAM learns from activated noise products as much as SGD. The activation patterns are further utilized to bound products and norms of the weight, signal and noise, which characterize the final test errors.

 Our task is reduced to bounding  $\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right)\leq 0)$ , then use previous error bounds for SGD in App. D.1 for  $\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{F}}(yf(\mathbf{W}^{(t)},\mathbf{x})>0)$ . The inner product with j=y can be bounded as

$$\left\langle \mathbf{w}_{y,r}^{(t)}, y\varphi \right\rangle = \left\langle \mathbf{w}_{y,r}^{(0)}, y\varphi \right\rangle + \kappa_{y,r}^{(t)} + \frac{1}{(P-1)} \sum_{i=1}^{n} \overline{\zeta}_{y,r,i}^{(t)} \cdot \|\boldsymbol{\xi}_{i}\|_{2}^{-2} \cdot \langle \boldsymbol{\xi}_{i}, y\varphi \rangle$$

$$+ \frac{1}{(P-1)} \sum_{i=1}^{n} \underline{\zeta}_{y,r,i}^{(t)} \cdot \|\boldsymbol{\xi}_{i}\|_{2}^{-2} \cdot \langle \boldsymbol{\xi}_{i}, y\varphi \rangle$$

$$\geq \left\langle \mathbf{w}_{y,r}^{(0)}, y\varphi \right\rangle + \kappa_{y,r}^{(t)}$$

$$- \frac{\sqrt{2 \log(6n/\delta)}}{P-1} \cdot \sigma_{p} \|\varphi\|_{2} \cdot (\sigma_{p}^{2}d/2)^{-1} \left[ \sum_{i=1}^{n} \overline{\zeta}_{y,r,i}^{(t)} + \sum_{i=1}^{n} \left| \underline{\zeta}_{y,r,i}^{(t)} \right| \right]$$

$$= \left\langle \mathbf{w}_{y,r}^{(0)}, y\varphi \right\rangle + \kappa_{y,r}^{(t)} - \Theta \left( \sqrt{\log(n/\delta)} \cdot (P\sigma_{p}d)^{-1} \|\varphi\|_{2} \right) \cdot \Theta \left( \text{SNR}^{-2} \right) \cdot \kappa_{y,r}^{(t)}$$

$$= \left\langle \mathbf{w}_{y,r}^{(0)}, y\varphi \right\rangle + \left[ 1 - \Theta \left( \sqrt{\log(n/\delta)} \cdot P\sigma_{p}/\|\varphi\|_{2} \right) \right] \kappa_{y,r}^{(t)}$$

$$= \left\langle \mathbf{w}_{y,r}^{(0)}, y\varphi \right\rangle + \Theta \left( \kappa_{y,r}^{(t)} \right) = \Theta(1),$$

$$(45)$$

where the inequality is by Lemma D.3; the second equality is obtained by plugging in the coefficient orders we summarized; the third equality is by SNR =  $\|\varphi\|_2/(P\sigma_p\sqrt{d})$ ; the fourth equality is by  $\|\varphi\|_2^2 \ge C \cdot P^2\sigma_p^2 \log(n/\delta)$  in Assumption D.1 for sufficiently large constant C; the last equality is by Lemma D.7 of Chen et al. (2023). We similarly have  $\langle \mathbf{w}_{y,r}^{(t)}, y\varphi \rangle = -\Theta(1) < 0$ .

Denote  $g(\xi)$  as  $\sum_{r} \sigma(\langle \mathbf{w}_{-y,r}^{(t)}, \xi \rangle)$ . The results for noise learning from SGD in App. D.1 still apply:

$$|g(\boldsymbol{\xi}) - g(\boldsymbol{\xi}')| \leq \sum_{r=1}^{m} \left\| \mathbf{w}_{-y,r}^{(t)} \right\|_{2} \cdot \|\boldsymbol{\xi} - \boldsymbol{\xi}'\|_{2};$$

$$\mathbb{E}g(\boldsymbol{\xi}) = \frac{\sigma_{p}}{\sqrt{2\pi}} \sum_{r=1}^{m} \left\| \mathbf{w}_{-y,r}^{(t)} \right\|_{2};$$

$$\left\| \sum_{i=1}^{n} \zeta_{j,r,i}^{(t)} \cdot \|\boldsymbol{\xi}_{i}\|_{2}^{-2} \cdot \boldsymbol{\xi}_{i} \right\|_{2}^{2} \leq \Theta\left(\sigma_{p}^{-2} d^{-1} n^{-1}\right) \left(\sum_{i=1}^{n} \overline{\zeta}_{j,r,i}^{(t)}\right)^{2}.$$
(46)

We can thus upper bound the 2-norm of  $\mathbf{w}_{i,r}^{(t)}$  as:

$$\left\| \mathbf{w}_{j,r}^{(t)} \right\|_{2} \leq \left\| \mathbf{w}_{j,r}^{(0)} \right\|_{2} + \kappa_{j,r}^{(t)} \cdot \left\| \varphi \right\|_{2}^{-1} + \frac{1}{P-1} \left\| \sum_{i=1}^{n} \zeta_{j,r,i}^{(t)} \cdot \left\| \boldsymbol{\xi}_{i} \right\|_{2}^{-2} \cdot \boldsymbol{\xi}_{i} \right\|_{2}$$

$$\leq \left\| \mathbf{w}_{j,r}^{(0)} \right\|_{2} + \kappa_{j,r}^{(t)} \cdot \left\| \varphi \right\|_{2}^{-1} + \Theta \left( P^{-1} \sigma_{p}^{-1} d^{-1/2} n^{-1/2} \right) \cdot \sum_{i=1}^{n} \overline{\zeta}_{j,r,i}^{(t)}$$

$$= \Theta(\sigma_{0} \sqrt{d}) + \Theta \left( P^{-1} \sigma_{p}^{-1} d^{-1/2} n^{-1/2} \right) \cdot \sum_{i=1}^{n} \overline{\zeta}_{j,r,i}^{(t)},$$

$$(47)$$

based on  $\mathrm{SNR} = \|\varphi\|_2/(P\sigma_p\sqrt{d})$  and  $\sum_{i=1}^n\overline{\zeta}_{j,r,i}^{(t)}/\kappa_{j,r}^{(t)} = \Theta\left(\mathrm{SNR}^{-2}\right)$ , and the condition for d in Assumption D.1, and also  $\left\|\mathbf{w}_{j,r}^{(0)}\right\|_2 = \Theta\left(\sigma_0\sqrt{d}\right)$  based on Lemma D.7 of Chen et al. (2023). Then

we have

$$\frac{\sum_{r} \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, y\varphi\right\rangle\right)}{(P-1)\sigma_{p} \sum_{r=1}^{m} \left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}} \geq \frac{\Theta(1)}{\Theta\left(\sigma_{0}\sqrt{d}\right) + \Theta\left(P^{-1}\sigma_{p}^{-1}d^{-1/2}n^{-1/2}\right) \cdot \sum_{i=1}^{n} \overline{\zeta}_{j,r,i}^{(t)}}$$

$$\geq \frac{\Theta(1)}{\Theta\left(\sigma_{0}\sqrt{d}\right) + O\left(P^{-1}\sigma_{p}^{-1}d^{-1/2}n^{1/2}\alpha\right)}$$

$$\geq \min\left\{\Omega\left(\sigma_{0}^{-1}d^{-1/2}\right), \Omega\left(P\sigma_{p}d^{1/2}n^{-1/2}\alpha^{-1}\right)\right\}$$

$$\geq 1$$

$$\Rightarrow \sum_{r} \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, y\varphi\right\rangle\right) - \frac{(P-1)\sigma_{p}}{\sqrt{2\pi}} \sum_{r=1}^{m} \left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2} > 0.$$
(48)

Upper bound. Now plug in previous results to obtain

$$\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right)\leq0\right)\leq\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left((P-1)\sum_{r}\sigma\left(\left\langle\mathbf{w}_{-y,r}^{(t)},\xi\right\rangle\right)\geq\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\varphi\right\rangle\right)\right)$$

$$=\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(g(\boldsymbol{\xi})-\mathbb{E}g(\boldsymbol{\xi})\geq1/(P-1)\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\varphi\right\rangle\right)-\frac{\sigma_{p}}{\sqrt{2\pi}}\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}\right)$$

$$\leq\exp\left[-\frac{c\left(1/(P-1)\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\varphi\right\rangle\right)-\left(\sigma_{p}/\sqrt{2\pi}\right)\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}\right)^{2}}{\sigma_{p}^{2}\left(\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}\right)^{2}}\right]$$

$$=\exp\left[-c\left(\frac{\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\varphi\right\rangle\right)}{(P-1)\sigma_{p}\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}}-1/\sqrt{2\pi}\right)^{2}\right]$$

$$\leq\exp(c/2\pi)\exp\left(-0.5c\left(\frac{\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\varphi\right\rangle\right)}{(P-1)\sigma_{p}\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}}\right)^{2}\right).$$

$$(49)$$

The second inequality is by Eq. 48 and plugging  $||g||_{\text{Lip}} \leq \sum_{r=1}^{m} ||\mathbf{w}_{-y,r}^{(t)}||_2$  into Eq. 15, the third inequality is because  $(s-t)^2 \geq s^2/2 - t^2, \forall s, t \geq 0$ . And we can obtain

$$\mathbb{P}_{(\mathbf{x},y)\sim\mathcal{R}}\left(yf\left(\mathbf{W}^{(t)},\mathbf{x}\right)\leq 0\right) \leq \exp(c/2\pi)\exp\left(-0.5c\left(\frac{\sum_{r}\sigma\left(\left\langle\mathbf{w}_{y,r}^{(t)},y\varphi\right\rangle\right)}{(P-1)\sigma_{p}\sum_{r=1}^{m}\left\|\mathbf{w}_{-y,r}^{(t)}\right\|_{2}}\right)^{2}\right) \\
\leq \exp\left(\frac{c}{2\pi}-C\min\left\{\sigma_{0}^{-2}d^{-1},P\sigma_{p}^{2}dn^{-1}\alpha^{-2}\right\}\right) \\
\leq \exp\left(-0.5C\min\left\{\sigma_{0}^{-2}d^{-1},P\sigma_{p}^{2}dn^{-1}\alpha^{-2}\right\}\right) = \epsilon,$$
(50)

where C=O(1), the last inequality holds since  $\sigma_0^2\leq 0.5Cd^{-1}\log(1/\epsilon)$  and  $d\geq 2C^{-1}P^{-1}\sigma_p^{-2}n\alpha^2\log(1/\epsilon)$ . Now we upper bound the test error  $\mathcal{L}^{\text{test}}(\mathbf{W}^{T_2},\mathcal{D})$ . Depending on the strength of the unified signal vector  $\varphi$ , the unlearning of  $\mathcal{F}$  can exhibit either benign or harmful overfitting following SGD's characterization, dividing error bounds into two cases:

1. If  $\|\varphi\|_2 \ge C_1 d^{1/4} n^{-1/4} P \sigma_p$ , we have benign overfitting on both  $\mathcal{R}$  and  $\mathcal{F}$ . Thus,

$$\mathcal{L}^{\text{test}}(\mathbf{W}^{T_2}, \mathcal{D}) = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left( y \neq \text{sign} \left( f\left( \mathbf{W}^{T_2}, \mathbf{x} \right) \right) \right)$$

$$= \beta \cdot \underbrace{\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{R}} \left( y f\left( \mathbf{W}^{T_2}, \mathbf{x} \right) \leq 0 \right)}_{\leq \epsilon_{\mathcal{R}}} + (1 - \beta) \cdot \left( 1 - \underbrace{\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{F}} \left( y f\left( \mathbf{W}^{T_2}, \mathbf{x} \right) > 0 \right)}_{\leq \epsilon_{\mathcal{F}}} \right)$$

$$\implies \lim_{\beta \to 1} \mathcal{L}^{\text{test}}(\mathbf{W}^{T_2}, \mathcal{D}) \leq \epsilon_{\mathcal{R}} = \epsilon.$$

(51)

As  $\beta \to 1$ ,  $|\mathcal{F}|/n$  decreases so the model can better maintain its performance; as  $\beta \to 0.5$ ,  $|\mathcal{F}|/n$  increases and more samples are to be unlearned, making the model performance reduce to a minimally PAC-learnable guarantee. Hence, when  $\beta \to 0.5$ , we have  $\lim_{\beta \to 0.5} \mathcal{L}^{\text{test}}(\mathbf{W}^{T_2}, \mathcal{D}) \leq 0.5 + 0.5\epsilon_{\mathcal{R}} - 0.5\epsilon_{\mathcal{F}} = \epsilon$ .

2. If  $\Omega(1) \leq \|\varphi\|_2 \leq C_1 d^{1/4} n^{-1/4} P \sigma_p$ , we have benign overfitting on  $\mathcal{R}$  and harmful overfitting on  $\mathcal{F}$ . Thus,

$$\mathcal{L}^{\text{test}}(\mathbf{W}^{T_2}, \mathcal{D}) = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left( y \neq \text{sign} \left( f\left(\mathbf{W}^{(t)}, \mathbf{x}\right) \right) \right)$$

$$= \beta \cdot \underbrace{\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{R}} \left( y f\left(\mathbf{W}^{(t)}, \mathbf{x}\right) \leq 0 \right)}_{\leq \epsilon_{\mathcal{R}}} + (1 - \beta) \cdot \left( 1 - \underbrace{\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{F}} \left( y f\left(\mathbf{W}^{(t)}, \mathbf{x}\right) > 0 \right)}_{\geq 0.1} \right)$$

$$\implies \lim_{\beta \to 1} \mathcal{L}^{\text{test}}(\mathbf{W}^{T_2}, \mathcal{D}) \leq \epsilon_{\mathcal{R}} = \epsilon.$$
(52)

Similarly, we have  $\lim_{\beta \to 0.5} \mathcal{L}^{\text{test}}(\mathbf{W}^{T_2}, \mathcal{D}) \leq 0.5 \epsilon_{\mathcal{R}} + 0.45 = \epsilon$ .

**Remark D.6** ( $\beta$ -dependence of the  $\epsilon$ -bound). The overall test error

$$\mathcal{L}^{\textit{test}}(\mathbf{W}^{T_2}, \mathcal{D}) = \beta \cdot \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{R}}\left(yf\left(\mathbf{W}^{(t)}, \mathbf{x}\right) \leq 0\right) + (1 - \beta) \cdot \left(1 - \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{F}}\left(yf\left(\mathbf{W}^{(t)}, \mathbf{x}\right) > 0\right)\right)$$

can be considered as an affine function of the mixing factor  $\beta$ , and so its achievable range runs from the best-case retain error  $\epsilon_{\mathcal{R}}$  (as  $\beta \to 1$ ) up to asymptotically 0.5 (as  $\beta \to 0.5$ )—the trivial PAC-learnability threshold. Concretely, by choosing  $\beta$  sufficiently close to 1, one drives  $\mathcal{L}^{test}(\mathbf{W}^{T_2}, \mathcal{D})$  arbitrarily close to the small "benign" error level  $\epsilon$ , whereas if  $\beta$  remains near 0.5 then  $\mathcal{L}^{test}(\mathbf{W}^{T_2}, \mathcal{D})$  can approach 0.5, the worst-case "minimally learnable" error. Thus, all our bounds interpolate smoothly between these two extremes via the single parameter  $\beta$ , and we report the most informative bounds in Theorem 3.2 and Theorem 3.3.

## D.3 PROOF TO COROLLARY 3.3.1

Recall the update rule for  $\kappa_{j,r}$ . For each epoch, the interference between retain and forget signals can be measured as

$$\sum_{b}^{|\mathcal{R}|/B} \alpha \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{R}}} \ell_{i}^{\prime(t,b)} \sigma'(\langle \mathbf{w}_{j,r}^{(t,b)}, y_{i} \boldsymbol{\varphi} \rangle) - \sum_{b}^{|\mathcal{F}|/B} (1 - \alpha) \frac{|\mathcal{R}|}{|\mathcal{F}|} \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{F}}} \ell_{i}^{\prime(t,b)} \sigma'(\langle \mathbf{w}_{j,r}^{(t,b)}, y_{i} \boldsymbol{\varphi} \rangle). \tag{53}$$

Similar to Lemma 3.1, the expected gradient values between retain and forget samples should not differ. Since we cycle the forget set to synchronously train with the retain set, updates from  $\mathcal{F}$  has been scaled up by  $\frac{|\mathcal{R}|}{|\mathcal{T}|}$ . Hence,

$$\mathbb{E}\left[\sum_{b}^{|\mathcal{R}|/B} \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{R}}} \ell_{i}^{\prime(t,b)} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t,b)}, y_{i} \boldsymbol{\varphi} \rangle)\right] = \mathbb{E}\left[\sum_{b}^{|\mathcal{F}|/B} \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{F}}} \ell_{i}^{\prime(t,b)} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t,b)}, y_{i} \boldsymbol{\varphi} \rangle)\right]$$
(54)

Combining together, to expect  $\kappa_{j,r}$  to increase monotonically every epoch, we want

$$\mathbb{E}\left[\sum_{b}^{|\mathcal{R}|/B} \alpha \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{R}}} \ell_{i}^{\prime(t,b)} \sigma'(\langle \mathbf{w}_{j,r}^{(t,b)}, y_{i} \boldsymbol{\varphi} \rangle) - \sum_{b}^{|\mathcal{F}|/B} (1-\alpha) \frac{|\mathcal{R}|}{|\mathcal{F}|} \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{F}}} \ell_{i}^{\prime(t,b)} \sigma'(\langle \mathbf{w}_{j,r}^{(t,b)}, y_{i} \boldsymbol{\varphi} \rangle)\right] \geq 0$$

$$\implies \alpha - (1-\alpha) \frac{|\mathcal{R}|}{|\mathcal{F}|} \geq 0 \implies \alpha \geq \frac{|\mathcal{R}|}{|\mathcal{F}| + |\mathcal{R}|}.$$
(55)

## D.4 PROOF TO LEMMA 3.4

By Theorem 3.3, SAM turns off noise memorization prevention mechanism when fitting  $\mathcal{F}$ , which leads to the same requirement on signal strength as SGD. The only difference between SAM and SGD under NegGrad is the more effective learning on  $\mathcal{R}$ . From Eq. 7 we have the per-batch update of  $\kappa_{i,r}$  on  $\mathcal{R}$  as

$$\Delta \kappa_{j,r} = \frac{\eta \|\boldsymbol{\varphi}\|_2^2}{Bm} \alpha \sum_{i \in \mathcal{I}_{t,b}^{\mathcal{R}}} \ell_i^{\prime(t,b)} \sigma^{\prime}(\langle \mathbf{w}_{j,r}^{(t,b)}, y_i \boldsymbol{\varphi} \rangle). \tag{56}$$

Let g denote the batch-average magnitude of  $\ell_i^{\prime(t,b)}\sigma'(\langle \mathbf{w}_{j,r}^{(t,b)},y_i\boldsymbol{\varphi}\rangle)$  for convenience. We can then express per-epoch  $\kappa$  update as

$$\Delta_{\operatorname{epoch}} \kappa_{j,r} = \frac{\eta \|\varphi\|_2^2}{m} \alpha |\mathcal{R}| g. \tag{57}$$

Now, consider achieving benign overfitting on  $\mathcal{R}$  only, where SGD requires  $\|\varphi\|_2 = \Omega(d^{1/4}|\mathcal{R}|^{-1/4}P\sigma_p)$  while SAM only requires  $\|\varphi\|_2 = \Omega(1)$ . That being said, given a fixed universal  $\varphi$  for  $\mathcal{D}$  and a choice of  $\alpha$ , we have SAM learning the retain signals faster than SGD:

$$\frac{\Delta_{\operatorname{epoch}} \kappa_{j,r}^{\operatorname{SAM}}}{\Delta_{\operatorname{epoch}} \kappa_{j,r}^{\operatorname{SGD}}} = \Theta(d^{1/2} |\mathcal{R}|^{-1/2} P^2 \sigma_p^2) = \Theta(\|\varphi\|_2^2). \tag{58}$$

Hence, in order to achieve the same signal learning performance as SAM on  $\mathcal{R}$ , SGD needs to scale up  $\alpha^{\text{SGD}}$ . Thus,

$$\frac{\alpha^{\text{SGD}}}{\alpha^{\text{SAM}}} = \Theta(d^{1/2}|\mathcal{R}|^{-1/2}P^2\sigma_p^2) = \Theta(\|\varphi\|_2^2), \text{ or } \alpha^{\text{SGD}} - \alpha^{\text{SAM}} = \Theta(\|\varphi\|_2^2). \tag{59}$$

In general, since  $|\mathcal{R}| = \Theta(n)$ , we can characterize the gap between  $\alpha^{\text{SGD}}$  and  $\alpha^{\text{SAM}}$  by  $O(\sqrt{d/n})$ .

## E IMPLEMENTATION DETAILS

#### E.1 EXPERIMENT SETUP

We conduct major experiments on CIFAR-100 (Krizhevsky et al., 2009) and ImageNet-1K (Russakovsky et al., 2015) using ResNet-50 (He et al., 2016). We adopt pre-computed memorization scores for these two datasets from Feldman & Zhang (2020) to generate  $\mathcal F$  of different memorization levels with  $|\mathcal F| \approx 5\% |\mathcal S|$ . We have  $|\mathcal F| = 3000$  for CIFAR-100 and  $|\mathcal F| = 60000$  for ImageNet. We sample high-memorization forget set  $\mathcal F_{high}$  by choosing  $|\mathcal F|$  samples of highest memorization scores from  $\mathcal S$ ,  $\mathcal F_{low}$  by choosing  $|\mathcal F|$  samples of lowest memorization scores, and  $\mathcal F_{mid}$  by choosing  $|\mathcal F|$  samples whose memorization scores are closest to 0.5. We also run experiments with randomly sampled  $\mathcal F_{rand}$  on Tiny-ImageNet and CIFAR-10 in App. G. We use RandomResizedCrop and RandomHorizontalFlip as train transforms.

**Pretraining and retraining.** We pretrain on  $\mathcal S$  and retrain on  $\mathcal R$  with the same settings. For CIFAR-100, we train for  $T_1=200$  epochs, use batch size 256, learning rate  $\eta_0=0.1$  with cosine annealing, SGD with momentum 0.9 and weight decay  $5\times 10^{-4}$ . For ImageNet, we train for  $T_1=150$  epochs, use batch size 512, learning rate  $\eta_0=0.25$  with cosine annealing and 5 warm-up epochs, SGD with momentum 0.9 and weight decay  $2\times 10^{-5}$ . For CIFAR-10, we train ResNet-18 for  $T_1=50$  epochs, use batch size 256, learning rate  $\eta_0=0.1$  with cosine annealing, SGD with momentum 0.9 and weight decay  $5\times 10^{-4}$ . We summarize the settings, test performance of different pretrained models, as well as accuracies of retrain models in Tab. 6.

Table 6: Differed settings of pretrained models and their test accuracies using different  $\mathcal{A}$  (top), as well as performance of retrained models w.r.t different  $\mathcal{F}$  (bottom) for CIFAR-100 and ImageNet-1K.

Dataset, Model	lr+warmup	Batch $B$	Epoch T	W. Decay	SGD	ASAM 0.1	ASAM 1.0	SAM 0.1
CIFAR100, Res50	0.1+0	256	200	5e-4	77.23	76.0	78.05	77.85
ImageNet, Res50	0.25+5	512	150	2e-5	75.04	74.94	76.53	76.18

Retrain	]	High Men	1		Mid Men	1		Low Mem	1
Dataset, Model	Retain	Forget	Test	Retain	Forget	Test	Retain	Forget	Test
CIFAR100, Res50 ImageNet, Res50	1	3.3 13.828	74.96 74.826	99.981 97.388		74.14 74.832	99.956 96.671	100.0 99.858	75.81 75.018

Unlearning. We conduct all unlearning methods for  $T_2=10$  epochs with the same batch size and optimizer settings. For NegGrad and Sharp MinMax, we unlearn with constant learning rate 0.02. We use  $\alpha=0.99$  for CIFAR-100 and  $\alpha=0.989$  for ImageNet accounting for its slightly smaller  $|\mathcal{F}|/|\mathcal{S}|$  ratio. For model splitting, we empirically find that a small ratio for forget model benefits ImageNet such as 5%, while CIFAR-100 suits a larger ratio such as 30%. For both pretraining and unlearning, we wrap SGD with vanilla SAM (Foret et al., 2020) with  $\rho=0.1$ , and Adaptive SAM (ASAM) (Kwon et al., 2021) with  $\rho=[0.1,1.0]$ , while keep other hyper-parameters the same for fair comparison.

#### E.2 SHARP MINMAX IMPLEMENTATION

Inspired by SalUn (Fan et al., 2023), we split the model into two and update using two separate optimizer, SAM and shaprness maximization. We split the model by ranking the parameters that are important to the forget set  $\mathcal F$  based on the magnitude of the gradient of the parameters after one pass on  $\mathcal F$ , and choose the highest percentage where we have 5% for ImageNet and 30% for CIFAR-100. Unlike SalUn, which essentially performs RL unlearning on the selected parameters, we update both models using opposite optimization. SalUn also requires a larger part of the model to fine-tune with noisy, label flipped  $\mathcal F$ . When running Sharp MinMax and SalUn, we load the weight mask corresponding to the loaded pretrained model for model splitting. We have summarized our implementation for weight masking in Alg. 1, and Sharp MinMax in Alg. 2.

## Algorithm 1 WeightMask

```
Require: forget_loader, model, criterion, percent
 1: for all (name, param) in model parameters do
 2:
        gradients[name] \leftarrow zeros\_like(param)
 4: for all (image, target) in forget_loader do
 5:
        loss \leftarrow criterion(model(image), target)
 6:
        optimizer.zero_grad(); loss.backward()
 7:
        accumulate parameter gradients into gradients
 8: end for
 9: for all name in gradients do
10:
        gradients[name] \leftarrow |gradients[name]|
11: end for
12: all_vals \leftarrow cat ({flatten(v) | v \in gradients.values()})
13: cutoff \leftarrow quantile(all\_vals, percent)
                                                                                 \triangleright e.g., 0.1 = bottom 10%
14: return { name \mapsto (grad < cutoff) | (name, grad) \in gradients}
```

#### E.3 UNLEARNING SETUP FOR PREVIOUS WORK

We compare with state-of-the-art unlearning methods with optimized hyper-parameter settings. To our best knowledge, several previous methods are evaluated on ImageNet for the first time. We apply SGD and ASAM 1.0 on each  $\mathcal U$  and compare the performance between SGD and SAM. For L1-Sparse (Jia et al., 2023), we use unlearn lr= 0.02 and  $\alpha = 1 \times 10^{-4}$ . For SCRUB (Kurmanji et al.,

1621

1644 1645

1646

1647

1648

1649 1650

1651 1652

1655

1656

1657

1658 1659

1662

1663

1664

1665

1666

1668

1669

1671

1672

1673

# Algorithm 2 SharpMinMax

```
Require: x_retain, y_retain, x_forget, y_forget, model, criterion, mask, alpha, optimizer_retain, op-
1622
                timizer_forget
1623
            1: r \cdot loss1 \leftarrow \alpha \cdot criterion(model(x \cdot retain), y \cdot retain)
1624
            2: r\_loss1.backward()
1625
            3: optimizer_retain.first_step(zero_grad=True)
                                                                                                                 ▷ SAM first step
1626
            4: r loss2 \leftarrow \alpha \cdot criterion(model(x_retain), y_retain)
1627
            5: r\_loss2.backward()
1628
            6: for all (name, p) in model parameters do
                    if p.grad then
1629
            7:
            8:
                         p.\operatorname{grad} \leftarrow p.\operatorname{grad} \odot (1 - \operatorname{mask[name]})
                                                                                                        1630
            9:
                     end if
           10: end for
1632
           11: optimizer_retain.second_step(zero_grad=True)
                                                                                                                       ⊳ sharp min
1633
           12: f \rfloor loss1 \leftarrow -(1-\alpha) \cdot criterion(model(x\_forget), y\_forget)
1634
           13: f\_loss1.backward()
1635
           14: optimizer_forget.first_step(zero_grad=True)

⊳ SAM first step

           15: f \cdot loss2 \leftarrow -(1-\alpha) \cdot \operatorname{criterion}(\operatorname{model}(\mathbf{x} \cdot \mathbf{forget}), \mathbf{y} \cdot \mathbf{forget})
1637
           16: f loss 2.backward()
           17: for all (name, p) in model parameters do
1639
           18:
                    if p.grad then
           19:
                         p.\operatorname{grad} \leftarrow p.\operatorname{grad} \odot \operatorname{mask}[\operatorname{name}]

    □ update forget params only

1640
           20:
                    end if
1641
           21: end for
1642
           22: optimizer_forget.second_step(zero_grad=True)
                                                                                                                       ⊳ sharp max
1643
```

2023), we use unlearn lr= 0.004, msteps= 8, kd\_T= 4,  $\beta=0.01$ , and  $\gamma=0.99$ . For RL (Graves et al., 2021), we use unlearn lr= 0.06 on CIFAR-100 and 0.02 on ImageNet. For SalUn (Fan et al., 2023), we use the unlearn lr= 0.06, 50% weight to finetune on CIFAR-100, and unlearn lr= 0.04, 30% weight to finetune on ImageNet.

#### E.4 EVALUATION DETAILS

**Membership inference attack.** We adopted a MIA based evaluation from Jia et al. (2023). We train a binary classifier using the retain set  $\mathcal{R}$  and the test set  $\mathcal{D}_{test}$  to distinguish whether a data sample was involved in the training stage, based on the softmaxed outputs from the unlearned model. Then, we feed the forget set  $\mathcal{F}$  to the classifier to evaluate this unlearned model. We expect forget samples to be classified as "non-training" data, and we evaluate the unlearning effectiveness based on MIA correctness. A lower correctness (close to 0.5) indicates difficulty to distinguish and thus better unlearning. This evaluation examines an unlearned model from a privacy perspective.

Entanglement computation. We compute both entanglement scores based on normalized embeddings of retain and forget sets from the penultimate layer of the model. We compute pair-wise entanglement between each retain and forget embedding, either globally or within a class. For variance-based entanglement  $E_{\text{Var}}$ , we directly follow Zhao et al. (2024) for implementation, and then rescale the raw scores to [0,1] based on the value range across global and class-wise scores. For Wasserstein entanglement  $E_{W_p}$ , we randomly sample an equal number of embeddings from retain and forget embeddings and build two uniform proxy-distributions. We then use existing optimal transport library to compute the transport distance (cost), outputting entanglement scores as 1- distance. No clipping is needed as we observe all scores lie within [0,1].

## F DETAILED EMPIRICAL RESULTS

## F.1 STATISTICAL SIGNIFICANCE

We demonstrate the statistical significance of our main empirical results by running each unlearning experiment three times with different seeds. In Fig. 4 and Fig. 3, we report the 95% confidence intervals ( $\mu \pm 2\sigma$ ) of all unlearning methods on ImageNet and CIFAR-100, which correspond to

Tab. 1 and Tab. 3. Each single bar represents the mean over runs and has the mean ToW scores marked on top of its error bar plotted by  $\pm 2\sigma$ . We observe that SAM consistently improves all unlearning methods with more noticeable results on CIFAR-100. For "All methods" subplots, we highlight the largest improvement by applying SAM to each  $\mathcal{U}$ . On CIFAR-100, we observe a general larger variance of SGD based unlearning, especially for SCRUB. Despite that  $\mathcal{A}=$ SAM 0.1 seems to provide a weaker pretrained model, Adaptive SAM settings can improve unlearning performance more steadily with lower variance, which demonstrate that SAM unlearning is more robust. Tab. 7 also records the means and variances of the "All methods" subplots for ImageNet and CIFAR-100. These additional insights further strengthen our findings.

Table 7: Verifying statistical significance ( $\mu\pm\sigma$ ) of main experiments on ImageNet and CIFAR-100. Given various pretrained model with different  $\mathcal{A}$ , we observe that SAM consistently improve base unlearn methods  $\mathcal{U}$  with higher means across multiple seeds. Moreover, we observe generally more stable performance with SAM based on smaller variance on average.

ImageNet	1	RL	Sa	ılUn	N	lG	Mir	nMax
Method	SGD	ASAM 1.0	SGD	ASAM 1.0	SGD	ASAM 1.0	SGD	ASAM 1.0
$\mathcal{A}$ =SGD	82.9±0.3	83.9±0.2	70.6±0.1	71.0±0.1	83.5±0.3	84.8±0.0	80.2±0.1	87.9±0.0
A=ASAM 0.1	82.5±0.1	$83.8 \pm 0.1$	$70.7\pm0.1$	$71.1 \pm 0.1$	83.4±0.3	$84.7 \pm 0.1$	$79.7\pm0.2$	$87.5 \pm 0.1$
A=ASAM 1.0	83.2±0.4	$83.8 \pm 0.2$	71.1±0.0	$71.2 \pm 0.0$	$84.1 \pm 0.0$	$84.6 \pm 0.2$	80.1±0.2	$88.0 \pm 0.1$
A=SAM 0.1	82.9±0.2	$83.7 \pm 0.3$	71.2±0.0	$71.4 \pm 0.1$	83.6±0.1	$84.4 \pm 0.1$	79.9±0.1	$87.8 \pm 0.1$

CIFAR100	L1 S	Sparse	Sc	rub	I	RL	Sa	lUn	N	NG
Method	SGD	ASAM 1.0	SGD	ASAM 1.0	SGD	ASAM 1.0	SGD	ASAM 1.0	SGD	ASAM 1.0
A=SGD	62.1±1.4	$67.3 \pm 0.1$	56.5±14.1	$73.6 \pm 0.4$	74.2±1.0	$77.2 \pm 0.2$	76.1±1.5	$83.8 {\pm} 0.9$	82.8±1.1	84.0±0.9
A=ASAM 0.1	63.6±1.7	$69.3 \pm 0.6$	54.3±1.8	$79.3 \pm 0.8$	$72.1\pm0.9$	$75.8 \pm 1.3$	72.9±1.6	$82.5\pm0.4$	83.9±0.8	$85.5 \pm 0.6$
A=ASAM 1.0	64.2±0.7	$68.7 \pm 1.7$	58.4±10.5	$72.0\pm2.1$	75.7±1.5	$80.3\pm1.2$	79.0±0.3	$83.3 \pm 0.2$	80.2±0.5	$83.9 \pm 0.2$
A=SAM 0.1	64.9±1.3	$68.3 \pm 0.6$	41.1±1.7	$49.7 \pm 16.6$	74.2±0.7	$80.3 \pm 0.9$	79.4±1.0	$83.6 \pm 0.6$	71.3±1.8	$78.7 \pm 0.5$

## F.2 COMPLETE ACCURACIES

In Tab. 8, Tab. 9, and Tab. 10, we report complete results of retain, forget, and test accuracies for all unlearning experiments, which are used to compute ToW scores in Tab. 1 and Tab. 3. As we have mentioned in the main paper, we observe that SGD often achieves lower test accuracies, motivating us to rethink the overfitting under a sample-specific unlearning scheme.

Table 8: Detailed accuracies of NegGrad on ImageNet and CIFAR-100.

High Mem   Retain   Forget   Test   ToW   Retain   Town   Test   Tow   Retain   Town   Test   Tow   Retain   Town   Test   Tow   Retain   Town   Test																	
+SGD	ImageNet		A =	SGD			A = AS	AM 0.1			A = AS	AM 1.0			A = SA	AM 0.1	
AASAM   1.   89.487   26.407   72.08   78.52   88.640   24.77   70.988   78.366   89.767   26.542   72.236   78.762   89.816   27.422   72.238   78.083   AASAM   1.0   99.804   28.398   73.506   77.898   90.498   28.445   73.05   78.301   91.583   30.997   73.746   77.388   91.323   31.578   73.767   77.762	High Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
+ASAM 1.0	+SGD	88.766	25.148		78.764		24.1	70.878				71.772	78.522	89.158	26.488	71.91	78.03
+SAM 0.1   91.007   29.88   73.676   77.898   90.498   28.445   73.05   78.301   91.583   30.997   73.746   77.388   91.328   31.578   73.964   76.807      Mid Mem	+ASAM 0.1	89.487	26.407	72.08		88.640		70.988						89.816	27.422	72.328	78.083
Hid Mem   Retain   Forget   Test   ToW   Retain   Test   ToW   Retain   Forget   Test   ToW   Retain   Test   ToW   Retain   Forget   Test   ToW   Retain   Test   Test   ToW   Retain   Test   ToW   Retain   Test   Test   ToW   Retain   Test   T																	
+SGD	+SAM 0.1	91.007	29.88	73.676	77.898	90.498	28.445	73.05	78.301	91.583	30.997	73.746	77.388	91.328	31.578	73.964	76.807
+ASAM 0.1 89.56 \$8.502 72.154 84.113 89.276 57.698 71.576 84.077 90.087 59.08 72.378 84.267 89.945 59.263 72.482 84.062   +ASAM 1.0 90.969 61.998 73.544 83.389 91.064 62.023 73.434 83.358 91.427 62.757 73.82 83.326 91.505 63.078 74.046 83.284   +SAM 0.1 91.396 63.015 73.734 82.985 91.015 62.308 73.422 83.04 91.984 64.367 74.014 82.473 91.823 64.258 74.198 82.587    Low Mem Retain Forget Test ToW Retain Forget RASAM 0.1 88.251 99.643 72.198 89.188 88.296 99.635 72.044 89.098 89.293 99.7 72.658 90.579 88.553 99.69 72.728 88.839 + ASAM 0.1 89.903 99.818 73.844 92.174 89.704 99.808 73.69 91.843 90.432 99.79 73.896 92.772 90.042 99.813 74.166 92.617 + SAM 0.1 90.234 99.822 74.21 92.841 89.553 99.817 73.728 91.722 90.815 99.827 74.228 93.429 90.184 99.825 74.254 92.829    CIFAR100	Mid Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
+ASAM 1.0   90,969   61,998   73,544   83,389   91,064   62,023   73,424   83,388   91,427   62,757   73,82   83,326   91,505   63,078   74,004   83,284   45AM 0.1   91,986   82,587   74,004   82,473   91,823   64,258   74,198   82,587   74,004   74,0014																	
High Mem   Retain   Forget   Test   ToW   Retain   Forget	+ASAM 0.1	89.56	58.502	72.154	84.113	89.276	57.698	71.576	84.07	90.087	59.08	72.378	84.267	89.945	59.263	72.482	84.062
Compage   Retain   Forget   Test   ToW   Retain   Test   Test   ToW   Retain   Test	+ASAM 1.0	90.969	61.998	73.544	83.389	91.064	62.023	73.434	83.358				83.326	91.505	63.078	74.046	83.284
+SGD	+SAM 0.1	91.396	63.015	73.734	82.985	91.015	62.308	73.422	83.04	91.984	64.367	74.014	82.473	91.823	64.258	74.198	82.587
+ASAM 0.1   88.251   99.643   72.198   89.188   88.296   99.635   72.044   89.098   89.293   99.7   72.658   90.579   88.553   99.69   72.776   89.973   +ASAM 1.0   89.903   99.812   73.844   92.174   89.704   99.808   73.69   91.843   90.432   99.79   73.896   92.772   90.042   99.813   74.166   92.617   +SAM 0.1   90.234   99.822   74.21   92.841   89.553   99.817   73.728   91.722   90.815   99.827   74.228   93.429   90.184   99.825   74.224   92.829      CIFAR100	Low Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
+ASAM 0.1   99.938   99.818   73.844   92.174   89.704   99.808   73.69   91.843   90.432   99.79   73.896   92.772   90.042   99.813   74.166   92.617	+SGD	87.775	99.617	71.942	88.515	86.592	99.505	71.042	86.651	88.847	99.663	72.41	89.947	87.847	99.625		88.839
SAM 0.1   90.234   99.822   74.21   92.841   89.553   99.817   73.728   91.722   90.815   99.827   74.228   93.429   90.184   99.825   74.254   92.829	+ASAM 0.1	88.251	99.643	72.198	89.188	88.296	99.635	72.044	89.098	89.293	99.7	72.658	90.579	88.553	99.69	72.776	89.973
CIFAR100         A = SGD         A = ASAM 0.1         A = ASAM 1.0         A = SAM 0.1           High Mem         Retain         Forget         Test         ToW         A = ASAM 0.1	+ASAM 1.0	89.903				89.704	99.808		91.843					90.042			
High Mem   Retain   Forget   Test   ToW   Retain   Test   Test   ToW   Retain   Test	+SAM 0.1	90.234	99.822	74.21	92.841	89.553	99.817	73.728	91.722	90.815	99.827	74.228	93.429	90.184	99.825	74.254	92.829
+SGD         92.929         12.9         68.17         78.334         94.05         11.433         66.68         79.277         94.533         15.267         67.78         77.274         91.814         22.4         66.23         67.82           +ASAM 0.1         93.736         13.467         67.71         78.131         94.852         11.633         67.32         80.336         94.633         15.333         67.82         77.331         93.674         22.9         67.94         70.054           +SAM 0.1         98.758         15.433         69.98         80.806         96.907         13.167         69.03         82.196         96.893         17.7         69.85         78.731         96.376         22.9         67.94         70.054           +SAM 0.1         98.555         19         72.82         81.331         99.193         17.4         72.17         82.86         99.4         26.467         72.74         47.04         99.24         36.767         73.49         65.08           Mid Mem         Retain         Forget         Test         ToW         Retain         Forget         Test         ToW         Retain         Forget         Test         ToW           +SGD         93.162																	
+SGD         92.929         12.9         68.17         78.334         94.05         11.433         66.68         79.277         94.533         15.267         67.78         77.274         91.814         22.4         66.23         67.82           +ASAM 0.1         93.736         13.467         67.71         78.131         94.852         11.633         67.32         80.336         94.633         15.333         67.82         77.331         93.674         22.9         67.94         70.054           +ASAM 1.0         96.748         15.433         69.98         80.806         96.907         13.167         69.03         82.196         96.893         17.7         69.85         78.731         96.376         24.033         69.85         72.518           #SAM 0.1         98.552         19         72.82         81.331         99.193         17.4         72.17         82.86         99.4         26.467         72.74         47.04         99.24         36.767         73.49         65.08           Mid Mem         Retain         Forget         Test         ToW         Retain         Forget         Test         ToW         Retain         Forget         Test         ToW           +SGD         93.162 </th <th>CIFAR100</th> <th>1</th> <th>A =</th> <th>SGD</th> <th></th> <th>1</th> <th>A = AS</th> <th>AM 0.1</th> <th></th> <th>1</th> <th>A = AS</th> <th>AM 1.0</th> <th></th> <th>1</th> <th>A = SA</th> <th>AM 0.1</th> <th></th>	CIFAR100	1	A =	SGD		1	A = AS	AM 0.1		1	A = AS	AM 1.0		1	A = SA	AM 0.1	
+ASAM 0.1 93.736 13.467 67.71 78.131 94.852 11.633 67.32 80.336 94.633 15.333 67.82 77.331 93.674 22.9 67.94 70.054 +ASAM 1.0 96.748 15.433 69.98 80.806 96.907 13.167 69.03 82.196 96.893 17.7 69.85 78.731 96.376 24.033 69.85 72.518 +SAM 0.1 98.552 19 72.82 81.331 99.193 17.4 72.17 82.86 99.4 26.467 72.74 74.704 99.24 36.767 73.49 65.08 Mid Mem Retain Forget Test ToW Retain Forget Retain Retain Retain Retain Forget Test ToW Retain Forget Retain Retain Retain Forget Test ToW Retain Forget Test ToW Retain Ret		Retain			ToW	Retain			ToW	Retain			ToW	   Retain			ToW
+ASAM 0.1 96.748 15.433 69.98 80.806 96.907 13.167 69.03 82.196 96.893 17.7 69.85 78.731 96.376 24.033 69.85 72.518 45.00 19.8552 19 72.82 81.331 99.193 17.4 72.17 82.86 99.4 26.467 72.74 74.704 99.24 36.767 73.49 65.08     Mid Mem	High Mem		Forget	Test			Forget	Test			Forget	Test			Forget	Test	
+SAM 0.1         98.552         19         72.82         81.331         99.193         17.4         72.17         82.86         99.4         26.467         72.74         74.704         99.24         36.767         73.49         65.08           Mid Mem         Retain         Forget         Test         ToW         Retain         Forget         Test         ToW         Retain         Forget         Test         ToW           +SGID         93.162         60.3         66.15         83.335         95.024         58.433         66.95         86.454         95.519         69.2         67.3         78.59         93.714         72.233         66.95         74.158           +ASAM 0.1         94.055         62.633         66.97         82.846         95.005         58.133         66.85         87.539         95.524         68.133         66.75         79.074         93.838         72.367         66.95         74.158           +ASAM 0.1         98.938         80.133         72.18         75.059         99.007         76.133         70.87         77.94         99.448         85.1         72.59         70.898         99.169         90.033         72.9         60.89           Low Mem         Reta	High Mem +SGD	92.929	Forget 12.9	Test 68.17	78.334	94.05	Forget 11.433	Test 66.68	79.277	94.533	Forget 15.267	Test 67.78	77.274	91.814	Forget 22.4	Test 66.23	67.82
+SGD   93.162   60.3   66.15   83.335   95.024   58.433   65.96   86.454   95.519   69.2   67.3   78.59   93.714   72.233   66.91   74.145   +ASAM 0.1   94.055   62.633   66.97   82.846   95.005   58.133   66.85   87.539   95.524   68.133   66.75   79.074   93.838   72.367   66.95   74.158   +ASAM 0.1   98.781   69.533   69.81   81.465   97.16   65.4   68.43   84.391   97.919   72.7   69.58   79.264   97.257   76.2   69.8   75.653   +SAM 0.1   98.938   80.133   72.18   75.059   99.007   76.133   70.87   77.94   99.448   85.1   72.59   70.898   99.169   90.033   72.9   66.089    Low Mem   Retain   Forget   Test   ToW   Retain   Forget   Test   ToW   Retain   Forget   Test   ToW   Retain   Forget   Test   ToW   +SGD   91.086   97.767   65.67   83.718   95.312   98.267   67.18   88.637   93.117   98.5   66.17   85.443   85.307   96.93   62.63   76.374   +ASAM 0.1   92.736   97.767   67.3   86.78   94.676   98.5   67.8   76.761   94.298   97.967   67.27   88.039   86.902   96.9   62.92   78.087   +ASAM 0.0   92.824   97.8   67.53   87.052   96.267   99.1   68.94   90.502   97.883   99.533   70.59   93.249   93.517   98.7   67.35   86.759	+SGD +ASAM 0.1	92.929 93.736	Forget 12.9 13.467	Test 68.17 67.71	78.334 78.131	94.05 94.852	Forget 11.433 11.633	Test 66.68 67.32	79.277 80.336	94.533 94.633	Forget 15.267 15.333	Test 67.78 67.82	77.274 77.331	91.814 93.674	Forget 22.4 22.9	Test 66.23 67.94	67.82 70.054
+ASAM 0.1   94.055   62.633   66.97   82.846   95.005   58.133   66.85   87.539   95.524   68.133   66.75   79.074   93.838   72.367   66.95   74.158   +ASAM 0.1   98.938   80.133   72.18   75.059   99.007   76.133   70.87   77.94   99.448   85.1   72.59   70.898   99.169   90.033   72.9   66.089	+SGD +ASAM 0.1 +ASAM 1.0	92.929 93.736 96.748	Forget 12.9 13.467 15.433	Test 68.17 67.71 69.98	78.334 78.131 80.806	94.05 94.852 96.907	Forget 11.433 11.633 13.167	Test 66.68 67.32 69.03	79.277 80.336 82.196	94.533 94.633 96.893	Forget 15.267 15.333 17.7	Test 67.78 67.82 69.85	77.274 77.331 78.731	91.814 93.674 96.376	Forget 22.4 22.9 24.033	Test 66.23 67.94 69.85	67.82 70.054 72.518
+ASAM 0.1 96.781 69.533 69.81 81.465 97.16 65.4 68.43 84.391 97.919 72.7 69.58 79.264 97.257 76.2 69.8 75.653 +SAM 0.1 98.938 80.133 72.18 75.059 99.007 76.133 70.87 77.94 99.448 85.1 72.59 70.898 99.169 90.033 72.9 66.089 Low Mem Retain Forget Test ToW Retain Forget Test	+SGD +ASAM 0.1 +ASAM 1.0 +SAM 0.1	92.929 93.736 96.748 98.552	Forget 12.9 13.467 15.433 19	Test 68.17 67.71 69.98 72.82	78.334 78.131 80.806 81.331	94.05 94.852 96.907 99.193	Forget 11.433 11.633 13.167 17.4	Test 66.68 67.32 69.03 72.17	79.277 80.336 82.196 82.86	94.533 94.633 96.893 99.4	Forget 15.267 15.333 17.7 26.467	Test 67.78 67.82 69.85 72.74	77.274 77.331 78.731 74.704	91.814 93.674 96.376 99.24	Forget 22.4 22.9 24.033 36.767	Test 66.23 67.94 69.85 73.49	67.82 70.054 72.518 65.08
+SAM 0.1         98.938         80.133         72.18         75.059         99.007         76.133         70.87         77.94         99.448         85.1         72.59         70.898         99.169         90.033         72.9         66.089           Low Mem         Retain         Forget         Test         ToW         Retain         Forget         Test         ToW         Retain         Forget         Test         ToW           +SGD         91.086         97.767         65.67         83.718         95.312         98.267         67.18         88.637         93.117         98.5         66.17         85.443         85.307         96.933         62.93         78.037           +ASAM 0.1         92.736         97.767         67.3         86.78         94.676         98.5         67         87.671         94.298         97.967         67.22         88.090         96.902         96.92         78.087           +ASAM 0.1         92.824         97.8         66.25         89.02         96.267         99.1         68.94         90.502         97.863         97.967         67.22         88.092         96.9         62.92         78.087           +ASAM 1.0         92.824         97.8         67.	+SGD +ASAM 0.1 +ASAM 1.0 +SAM 0.1 Mid Mem	92.929 93.736 96.748 98.552 Retain	12.9 13.467 15.433 19 Forget	Test 68.17 67.71 69.98 72.82 Test	78.334 78.131 80.806 81.331 ToW	94.05 94.852 96.907 99.193 Retain	Forget 11.433 11.633 13.167 17.4 Forget	Test 66.68 67.32 69.03 72.17 Test	79.277 80.336 82.196 82.86 ToW	94.533 94.633 96.893 99.4 Retain	Forget 15.267 15.333 17.7 26.467 Forget	Test 67.78 67.82 69.85 72.74 Test	77.274 77.331 78.731 74.704 ToW	91.814 93.674 96.376 99.24 Retain	22.4 22.9 24.033 36.767 Forget	Test 66.23 67.94 69.85 73.49 Test	67.82 70.054 72.518 65.08
Low Mem         Retain         Forget         Test         ToW           4-SAM 0.1         92.736         97.767         67.3         86.782         98.267         98.5         67         87.671         94.298         97.967         67.27         88.039         86.902         96.96         62.92         78.087           4ASAM 0.1         92.824         97.8         67.53         87.052         96.269         99.1         68.94         90.502         97.883         99.533         70.59         32.249         93.51         98.7         67.35         86.759	+SGD +ASAM 0.1 +ASAM 0.1 +SAM 0.1 Mid Mem +SGD	92.929 93.736 96.748 98.552 Retain	Forget 12.9 13.467 15.433 19 Forget 60.3	Test 68.17 67.71 69.98 72.82 Test 66.15	78.334 78.131 80.806 81.331 ToW 83.335	94.05 94.852 96.907 99.193 Retain	Forget 11.433 11.633 13.167 17.4 Forget 58.433	Test 66.68 67.32 69.03 72.17 Test 65.96	79.277 80.336 82.196 82.86 ToW 86.454	94.533 94.633 96.893 99.4 Retain	Forget 15.267 15.333 17.7 26.467 Forget 69.2	Test 67.78 67.82 69.85 72.74 Test 67.3	77.274 77.331 78.731 74.704 ToW 78.59	91.814 93.674 96.376 99.24 Retain	Forget 22.4 22.9 24.033 36.767 Forget 72.233	Test 66.23 67.94 69.85 73.49 Test 66.91	67.82 70.054 72.518 65.08 ToW 74.145
+SGD   91.086   97.767   65.67   83.718   95.312   98.267   67.18   88.637   93.117   98.5   66.17   85.443   85.307   96.933   62.63   76.374   +ASAM 0.1   92.736   97.767   67.3   86.78   94.676   98.5   67   87.671   94.298   97.967   67.27   88.039   86.902   96.9   62.92   78.087   +ASAM 1.0   92.824   97.8   67.53   87.052   96.267   99.1   68.94   90.502   97.883   99.533   70.59   93.249   93.517   98.7   67.35   86.759	+SGD +ASAM 0.1 +ASAM 0.1 +SAM 0.1 Mid Mem +SGD +ASAM 0.1	92.929 93.736 96.748 98.552 Retain 93.162 94.055	Forget 12.9 13.467 15.433 19 Forget 60.3 62.633	Test 68.17 67.71 69.98 72.82 Test 66.15 66.97	78.334 78.131 80.806 81.331 ToW 83.335 82.846	94.05 94.852 96.907 99.193 Retain 95.024 95.005	Forget 11.433 11.633 13.167 17.4 Forget 58.433 58.133	Test 66.68 67.32 69.03 72.17 Test 65.96 66.85	79.277 80.336 82.196 82.86 ToW 86.454 87.539	94.533 94.633 96.893 99.4 Retain 95.519 95.524	Forget 15.267 15.333 17.7 26.467 Forget 69.2 68.133	Test 67.78 67.82 69.85 72.74 Test 67.3 66.75	77.274 77.331 78.731 74.704 ToW 78.59 79.074	91.814 93.674 96.376 99.24   Retain   93.714 93.838	Forget 22.4 22.9 24.033 36.767 Forget 72.233 72.367	Test 66.23 67.94 69.85 73.49 Test 66.91 66.95	67.82 70.054 72.518 65.08 ToW 74.145 74.158
+ASAM 0.1   92.736   97.767   67.3   86.78   94.676   98.5   67   87.671   94.298   97.967   67.27   88.039   86.902   96.9   62.92   78.087   +ASAM 1.0   92.824   97.8   67.53   87.052   96.267   99.1   68.94   90.502   97.883   99.533   70.59   93.249   93.517   98.7   67.35   86.759	+SGD +ASAM 0.1 +ASAM 1.0 +SAM 0.1 Mid Mem +SGD +ASAM 0.1 +ASAM 1.0	92.929 93.736 96.748 98.552 Retain 93.162 94.055 96.781	Forget 12.9 13.467 15.433 19 Forget 60.3 62.633 69.533	Test 68.17 67.71 69.98 72.82 Test 66.15 66.97 69.81	78.334 78.131 80.806 81.331 ToW 83.335 82.846 81.465	94.05 94.852 96.907 99.193 Retain 95.024 95.005 97.16	Forget 11.433 11.633 13.167 17.4 Forget 58.433 58.133 65.4	Test 66.68 67.32 69.03 72.17 Test 65.96 66.85 68.43	79.277 80.336 82.196 82.86 ToW 86.454 87.539 84.391	94.533 94.633 96.893 99.4 Retain 95.519 95.524 97.919	Forget 15.267 15.333 17.7 26.467 Forget 69.2 68.133 72.7	Test 67.78 67.82 69.85 72.74 Test 67.3 66.75 69.58	77.274 77.331 78.731 74.704 ToW 78.59 79.074 79.264	91.814 93.674 96.376 99.24 Retain 93.714 93.838 97.257	Forget 22.4 22.9 24.033 36.767 Forget 72.233 72.367 76.2	Test 66.23 67.94 69.85 73.49 Test 66.91 66.95 69.8	67.82 70.054 72.518 65.08 ToW 74.145 74.158 75.653
+ASAM 1.0   92.824   97.8   67.53   87.052   96.267   99.1   68.94   90.502   97.883   99.533   70.59   93.249   93.517   98.7   67.35   86.759	#Igh Mem +SGD +ASAM 0.1 +ASAM 1.0 +SAM 0.1  Mid Mem +SGD +ASAM 0.1 +ASAM 0.1 +ASAM 0.1	92.929 93.736 96.748 98.552 Retain 93.162 94.055 96.781 98.938	Forget 12.9 13.467 15.433 19 Forget 60.3 62.633 69.533 80.133	Test 68.17 67.71 69.98 72.82 Test 66.15 66.97 69.81 72.18	78.334 78.131 80.806 81.331 ToW 83.335 82.846 81.465 75.059	94.05 94.852 96.907 99.193 Retain 95.024 95.005 97.16 99.007	Forget 11.433 11.633 13.167 17.4 Forget 58.433 58.133 65.4 76.133	Test 66.68 67.32 69.03 72.17 Test 65.96 66.85 68.43 70.87	79.277 80.336 82.196 82.86 ToW 86.454 87.539 84.391 77.94	94.533 94.633 96.893 99.4 Retain 95.519 95.524 97.919 99.448	Forget 15.267 15.333 17.7 26.467 Forget 69.2 68.133 72.7 85.1	Test 67.78 67.82 69.85 72.74 Test 67.3 66.75 69.58 72.59	77.274 77.331 78.731 74.704 ToW 78.59 79.074 79.264 70.898	91.814 93.674 96.376 99.24 Retain 93.714 93.838 97.257 99.169	Forget  22.4 22.9 24.033 36.767  Forget 72.233 72.367 76.2 90.033	Test 66.23 67.94 69.85 73.49 Test 66.91 66.95 69.8 72.9	67.82 70.054 72.518 65.08 ToW 74.145 74.158 75.653 66.089
	#SGD +ASAM 0.1 +ASAM 1.0 +SAM 0.1 #Mid Mem +SGD +ASAM 0.1 +ASAM 1.0 +SAM 0.1 Low Mem	92.929 93.736 96.748 98.552 Retain 93.162 94.055 96.781 98.938 Retain	Forget 12.9 13.467 15.433 19 Forget 60.3 62.633 69.533 80.133 Forget	Test 68.17 67.71 69.98 72.82 Test 66.15 66.97 69.81 72.18 Test	78.334 78.131 80.806 81.331 ToW 83.335 82.846 81.465 75.059	94.05 94.852 96.907 99.193 Retain 95.024 95.005 97.16 99.007	Forget 11.433 11.633 13.167 17.4 Forget 58.433 58.133 65.4 76.133 Forget	Test 66.68 67.32 69.03 72.17 Test 65.96 66.85 68.43 70.87 Test	79.277 80.336 82.196 82.86 ToW 86.454 87.539 84.391 77.94 ToW	94.533 94.633 96.893 99.4   Retain   95.519 95.524 97.919 99.448   Retain	Forget 15.267 15.333 17.7 26.467 Forget 69.2 68.133 72.7 85.1 Forget	Test 67.78 67.82 69.85 72.74 Test 67.3 66.75 69.58 72.59	77.274 77.331 78.731 74.704 ToW 78.59 79.074 79.264 70.898 ToW	91.814 93.674 96.376 99.24   Retain   93.714 93.838 97.257 99.169   Retain	Forget  22.4  22.9  24.033  36.767  Forget  72.233  72.367  76.2  90.033  Forget	Test 66.23 67.94 69.85 73.49 Test 66.91 66.95 69.8 72.9 Test	67.82 70.054 72.518 65.08 ToW 74.145 74.158 75.653 66.089
	## High Mem  +SGD +ASAM 0.1 +ASAM 1.0 +SAM 0.1  Mid Mem +SGD +ASAM 0.1 +ASAM 1.0  Low Mem +SGD	92.929 93.736 96.748 98.552 Retain 93.162 94.055 96.781 98.938 Retain 91.086	Forget 12.9 13.467 15.433 19 Forget 60.3 62.633 69.533 80.133 Forget 97.767	Test 68.17 67.71 69.98 72.82 Test 66.15 66.97 69.81 72.18 Test	78.334 78.131 80.806 81.331 ToW 83.335 82.846 81.465 75.059 ToW 83.718	94.05 94.852 96.907 99.193 Retain 95.024 95.005 97.16 99.007 Retain 95.312	Forget 11.433 11.633 13.167 17.4 Forget 58.433 65.4 76.133 Forget 98.267	Test 66.68 67.32 69.03 72.17 Test 65.96 66.85 68.43 70.87 Test 67.18	79.277 80.336 82.196 82.86 ToW 86.454 87.539 84.391 77.94 ToW 88.637	94.533 94.633 96.893 99.4 Retain 95.519 95.524 97.919 99.448 Retain 93.117	Forget 15.267 15.333 17.7 26.467 Forget 69.2 68.133 72.7 85.1 Forget 98.5	Test 67.78 67.82 69.85 72.74 Test 67.3 66.75 69.58 72.59 Test 66.17	77.274 77.331 78.731 74.704 ToW 78.59 79.074 79.264 70.898 ToW 85.443	91.814 93.674 96.376 99.24 Retain 93.714 93.838 97.257 99.169 Retain 85.307	Forget 22.4 22.9 24.033 36.767 Forget 72.233 72.367 76.2 90.033 Forget 96.933	Test 66.23 67.94 69.85 73.49 Test 66.91 66.95 69.8 72.9 Test 62.63	67.82 70.054 72.518 65.08 ToW 74.145 74.158 75.653 66.089 ToW 76.374
	## High Mem  +SGD +ASAM 0.1 +ASAM 1.0 +SAM 0.1  ## Mid Mem +SGD +ASAM 0.1 +ASAM 0.1  Low Mem +SGD +ASAM 0.1	92.929 93.736 96.748 98.552 Retain 93.162 94.055 96.781 98.938 Retain 91.086 92.736	Forget 12.9 13.467 15.433 19 Forget 60.3 62.633 69.533 80.133 Forget 97.767	Test 68.17 67.71 69.98 72.82 Test 66.15 66.97 69.81 72.18 Test 65.67 67.3	78.334 78.131 80.806 81.331 ToW 83.335 82.846 81.465 75.059 ToW 83.718 86.78	94.05 94.852 96.907 99.193 Retain 95.024 95.005 97.16 99.007 Retain 95.312 94.676	Forget 11.433 11.633 13.167 17.4 Forget 58.433 58.133 65.4 76.133 Forget 98.267 98.5	Test 66.68 67.32 69.03 72.17 Test 65.96 66.85 68.43 70.87 Test 67.18	79.277 80.336 82.196 82.86 ToW 86.454 87.539 84.391 77.94 ToW 88.637 87.671	94.533   94.633   96.893   99.4   Retain   95.519   95.524   97.919   99.448   Retain   93.117   94.298	Forget 15.267 15.333 17.7 26.467 Forget 69.2 68.133 72.7 85.1 Forget 98.5 97.967	Test 67.78 67.82 69.85 72.74 Test 67.3 66.75 69.58 72.59 Test 66.17 67.27	77.274 77.331 78.731 74.704 ToW 78.59 79.074 79.264 70.898 ToW 85.443 88.039	91.814 93.674 96.376 99.24   Retain   93.714 93.838 97.257 99.169   Retain   85.307 86.902	Forget 22.4 22.9 24.033 36.767 Forget 72.233 72.367 76.2 90.033 Forget 96.933 96.9	Test 66.23 67.94 69.85 73.49 Test 66.91 66.95 69.8 72.9 Test 62.63 62.92	67.82 70.054 72.518 65.08 ToW 74.145 75.653 66.089 ToW 76.374 78.087

Table 9: Detailed accuracies of Sharp MinMax on ImageNet and CIFAR-100.

SGD																	
+SGD   87.513   29.79   71.408   73.357   86.802   28.42   70.692   73.418   88.411   31.423   72.016   73.103   87.879   30.953   71.964   73.0   73.488   73.48   73			A = S	SGD			A = ASA	M 0.1			A = AS.	AM 1.0			A = SA	M 0.1	
+ASAM 0.1   79.741   10.555   66.334   78.066   80.84185   11.222   66.894   79.077   73.491   8.203   61.802   70.148   80.16741   11.032   66.828   78.5   +ASAM 0.1   87.993   15.903   72.224   86.6658   87.748   15.605   71.638   86.166   88.633   72.452   86.915   88.435   17.083   72.498   86.2   +SAM 0.1   88.297   16.705   72.48   86.463   87.537   16.098   71.612   85.511   89.056   17.405   72.812   86.849   88.468   17.92   72.674   85.7      Mid Mem   Retain   Forget   Test   ToW   Retain   Forget   Test   ToW   Retain   Forget   Test   ToW   Retain   Forget   Test   ToW   +SGD   87.089   58.915   71.418   80.881   86.757   58.372   71.1   80.784   87.217   59.095   71.734   81.105   87.461   59.677   71.848   80.9   +ASAM 0.1   88.679   54.642   72.834   87.345   88.588   54.548   72.666   87.192   89.12   55.377   73.018   87.27   89.092   55.733   73.192   87.0   +SAM 0.1   89.141   56.215   73.268   86.755   88.642   55.303   72.74   86.635   89.492   56.813   73.47   86.722   89.758   57.657   73.792   86.4      Low Mem   Retain   Forget   Test   ToW	High Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
+ASAM 1.0   87.993   15.903   72.224   86.658   87.748   15.605   71.638   86.166   88.563   16.453   72.452   86.915   88.435   17.083   72.498   86.249   88.400   17.92   72.674   85.754   85.714   86.297   16.705   72.812   86.849   88.468   17.92   72.674   85.754   85	+SGD	87.513	29.79	71.408	73.357	86.802	28.42	70.692	73.418	88.411	31.423	72.016	73.103	87.879	30.953	71.964	73.052
Heating   Heat																	78.529
Mid Mem   Retain   Forget   Test   ToW   Retain   Forget   T																	86.272
+SGD   87.089   58.915   71.418   80.881   86.757   58.372   71.1   80.784   87.217   59.995   71.734   81.105   87.461   59.677   71.848   80.9482   80.9482   80.936   50.585   71.38   87.914   86.281   49.833   70.814   87.40   87.561   51.3   71.528   88.093   87.529   52.043   71.84   87.64   87.840   87.840   87.840   87.940   88.679   54.642   72.834   87.345   88.588   58.548   72.666   87.192   89.12   55.377   73.018   87.27   89.92   55.333   73.192   87.040   87.84	+SAM 0.1	88.297	16.705	72.48	86.463	87.537	16.098	71.612	85.511	89.056	17.405	72.812	86.849	88.468	17.92	72.674	85.712
+ASAM 0.1 86.936 50.585 71.38 87.914 86.281 49.833 70.814 87.40 87.561 51.3 71.528 88.039 87.529 52.043 71.84 87.64 +ASAM 1.0 88.679 54.642 72.834 87.345 88.588 54.548 72.666 87.192 89.12 55.377 73.018 87.27 89.092 55.733 73.192 87.04 +ASAM 1.0 89.141 56.215 73.268 86.755 88.642 55.303 72.74 86.635 89.492 56.817 73.018 87.27 89.092 55.733 73.792 86.4      Low Mem   Retain   Forget   Test   ToW   Ret	Mid Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
+ASAM 1.0   88.679   54.642   72.834   87.345   88.588   54.548   72.666   87.192   89.12   55.377   73.018   87.27   89.092   55.733   73.192   87.0   +SAM 0.1   89.141   56.215   73.268   86.755   88.642   55.303   72.74   86.635   89.492   56.813   73.47   86.722   89.758   57.657   73.792   86.4   +SGD   85.798   99.61   71.644   86.334   84.348   99.482   70.894   84.378   85.863   99.568   71.61   86.402   85.098   99.57   71.45   85.5   +ASAM 0.1   86.399   99.565   72.07   87.338   86.236   99.562   71.814   86.953   86.644   99.627   72.104   87.554   85.894   99.597   71.45   85.5   +ASAM 0.1   87.766   99.768   73.392   89.694   87.366   99.772   73.216   89.138   88.159   99.722   73.412   90.142   87.837   99.765   73.718   90.0   +SAM 0.1   87.836   99.777   73.666   90.005   87.745   99.76   73.58   89.852   88.706   99.783   73.94   91.111   87.974   99.792   73.752   90.2      CIFAR100   A = SGD																	80.913
+SAM 0.1   89.141   56.215   73.268   86.755   88.642   55.303   72.74   86.635   89.492   56.813   73.47   86.722   89.758   57.657   73.792   86.4    Low Mem   Retain   Forget   Test   ToW   Retain   Test   ToW   Re																	87.642
Low Mem   Retain   Forget   Test   ToW   Retain   Test																	87.076
+SGD	+SAM 0.1	89.141	56.215	73.268	86.755	88.642	55.303	72.74	86.635	89.492	56.813	73.47	86.722	89.758	57.657	73.792	86.486
+ASAM 0.1 86.399 99.565 72.07 87.338 86.236 99.562 71.814 86.953 86.644 99.627 72.104 87.554 85.894 99.593 71.898 86.6 +ASAM 1.0 87.766 99.768 73.392 89.694 87.366 99.772 73.216 89.138 81.59 99.722 73.412 90.142 87.837 99.765 73.718 90.0  CIFAR100   A = SGD	Low Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
+ASAM 1.0 87.766 99.768 73.392 89.694 87.366 99.772 73.216 89.138 88.159 99.722 73.412 90.142 87.837 99.765 73.718 90.0   +SAM 0.1 87.836 99.777 73.666 90.005 87.745 99.76 73.58 89.852 88.706 99.783 73.94 91.111 87.974 99.792 73.752 90.2    CIFAR100	+SGD	85.798				84.348	99.482	70.894	84.378					85.098	99.57	71.45	85.517
+SAM 0.1   87.836   99.777   73.666   90.005   87.745   99.76   73.58   89.852   88.706   99.783   73.94   91.111   87.974   99.792   73.752   90.2      CIFAR100   A = SGD   A = ASAM 0.1   A = ASAM 0.1   A = ASAM 1.0   A = ASAM 0.1	+ASAM 0.1	86.399	99.565	72.07	87.338	86.236	99.562	71.814	86.953	86.644	99.627	72.104	87.554	85.894	99.593	71.898	86.668
CIFAR100         A = SGD         A = ASAM 0.1         A = ASAM 1.0         A = ASAM 1.0         A = SAM 0.1           High Mem         Retain         Forget         Test         ToW           +SGD         92.298         20.8         67.86         70.767         95.098         22.167         68.42         72.137         92.564         25.4         66.35         65.925         87.195         25.233         63.77         60.4	+ASAM 1.0								89.138								90.064
High Mem         Retain         Forget         Test         ToW         Retain         Forget         Test         ToW         Retain         Forget         Test         ToW           +SGD         92.298         20.8         67.86         70.767         95.098         22.167         68.42         72.137         92.564         25.4         66.35         65.925         87.195         25.233         63.77         60.4	+SAM 0.1	87.836	99.777	73.666	90.005	87.745	99.76	73.58	89.852	88.706	99.783	73.94	91.111	87.974	99.792	73.752	90.207
+SGD   92.298   20.8   67.86   70.767   95.098   22.167   68.42   72.137   92.564   25.4   66.35   65.925   87.195   25.233   63.77   60.4	CIFAR100	4 =SGD			1	A = ASAM 0.1				A = ASAM 1.0				A = SAM 0.1			
+SGD   92.298   20.8   67.86   70.767   95.098   22.167   68.42   72.137   92.564   25.4   66.35   65.925   87.195   25.233   63.77   60.4	High Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
																	60.478
																	72.897
+ASAM 1.0   92.121   6.467   67.15   82.27   88.976   5.067   63.68   77.576   93.895   6.567   67.98   84.521   90.448   10.7   65.71   76.0																	76.037
																	85.195
Mid Mem   Retain Forget Test ToW   Retain Forg	Mid Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
	+SGD	91.433		65.79	76.692	91.633	63.367	64.39	77.864	92.11	69.4	65.96	74.526	85.714		62.55	71.931
	+ASAM 0.1	91.16	42.7	65.88	96.027	91.4	40.233	64.11	96.451		51.2	66.74	93.786	88.074	55.867	63.61	80.104
																	83.633
	+SAM 0.1	97.433	60.867	70.73	90.96	97.874	55.033	69.39	95.543		64.333	70.62	88.646	98.824	76.433	71.84	78.286
Low Mem   Retain Forget Test ToW   Retain Forg	Low Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
	+SGD	89.579	97.6	66.09	82.853	89.781	97.1	64.36	81.847	88.605	97.833	64.28	80.127	81.488	94.467	61.63	73.843
+SGD   89.579   97.6   66.09   82.853   89.781   97.1   64.36   81.847   88.605   97.833   64.28   80.127   81.488   94.467   61.63   73.8		89.026	95.067	65.12		89.874	96.033	64.47		93.748			87.151	92,967	97.433		86.659
																	77.461
+ASAM 0.1 89.026 95.067 65.12 83.473 89.874 96.033 64.47 82.883 93.748 97.167 66.17 87.151 92.967 97.433 66.65 86.6	+SAM 0.1	96.129	98.033	70.13	92,494	96.829	98.7	69.06	91.508	97.624	98.567	69.85	93.163	96,652	99.033	68.98	90.963

## G ADDITIONAL EXPERIMENTS

We provide additional experiments on CIFAR-10 and Tiny-ImageNet using randomly sampled forget set  $\mathcal{F}_{rand}$ . To diversify our experiment settings, we use ResNet-34 with ImageNet-pretrained weights for our learning and unlearning on Tiny-ImageNet. Similar to our main setup, we pretrain and retrain using the same settings, and we have summarized basic settings and baseline performance in Tab. 11. Since Tiny-ImageNet has 100K samples, we set  $|\mathcal{F}_{rand}| = 6000$  for Tiny-ImageNet. Tab. 12 records detailed accuracies and ToW scores of various unlearning and pretraining settings.

#### G.1 CIFAR-10

We summarize detailed unlearning settings on CIFAR-10. For L1-Sparse, we use unlearn lr= 0.02 and  $\alpha=1\times10^{-4}$ . For SCRUB, we use unlearn lr= 0.004, msteps= 8, kd\_T= 3.5,  $\beta=0.01$ , and  $\gamma=0.99$ . For RL and SalUn, we use unlearn lr= 0.08, and use 50% model parameters for SalUn. For NegGrad and Sharp MinMax, we use unlearn lr= 0.02 and  $\alpha=0.99$ , and use 30% model parameters for unlearning on  $\mathcal F$  and the rest for learning on  $\mathcal R$ .

From the results in Tab. 11, we observe consistent improvement by using SAM except only two cases for RL and SalUn with  $\mathcal{A}=SGD$ . Surprisingly, Sharp MinMax is not the best algorithm on CIFAR-10. By the nature of its design to overfit to forget signals deliberately, we hypothesize that this approach might be aggressive for small-scale unlearning. We again observe SCRUB to be an unstable algorithm which collapses when unlearning with SGD given  $\mathcal{A}=SAM0.1$ , while SAM helps reduce variance and stabilizes SCRUB unlearning given various pretrained models.

## G.2 TINY-IMAGENET

We summarize detailed unlearning settings on Tiny-ImageNet. For L1-Sparse, we use unlearn lr= 0.002 and  $\alpha=1\times10^{-4}$ . For SCRUB, we use unlearn lr= 0.002, msteps= 8, kd\_T= 3.5,  $\beta=0.01$ , and  $\gamma=0.99$ . For RL and SalUn, we use unlearn lr= 0.015, and use 30% model parameters for SalUn. For NegGrad and Sharp MinMax, we use unlearn lr= 0.005 and  $\alpha=0.99$ , and use 10% model parameters for unlearning on  $\mathcal F$  and the rest for learning on  $\mathcal R$ .

From the results in Tab. 11, we observe consistent improvement by using SAM except few cases. SCRUB performs more steadily than on CIFAR-10. While RL and SalUn perform well on other datasets, they do not appear to be effective on Tiny-ImageNet.

Table 10: Detailed accuracies of previous methods on ImageNet and CIFAR-100.

ImageNet		A =	SGD			A = AS	AM 0.1			A = AS	AM 1.0		A = SAM 0.1			
High Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
RL	88.536	29.857	72.02	74.598	88.663	29.622	71.95	74.857	88.975	30.59	72.04	74.317	89.429	31.74	72.572	74.055
+ASAM 1.0	90.874	33.395	74.234	74.951	90.615	32.668	73.972	75.221	91.14	34.745	74.298	73.95	91.155	35.332	74.522	73.579
SalUn	93.248	67.118	75.04	44.981	93.016	65.807	74.976	46.104	93.124	66.372	75.418	45.814	92.911	66.333	75.982	46.006
+ASAM 1.0	93.123	66.217	75.496	45.998	92.963	65.058	75.28	46.938	93.134	66.472	75.712	45.856	92.855	66.032	76.172	46.358
Mid Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
RL	88.785	54.653	71.916	86.617	88.067	53.387	71.258	86.462	89.754	56.17	72.634	86.813	88.609	54.608	72.168	86.715
+ASAM 1.0	90.597	59.53	73.836	85.581	90.457	59.337	73.654	85.473	90.993	60.35	74.078	85.393	90.902	60.402	74.348	85.494
SalUn	93.174	77.258	74.816	71.839	93.072	77.222	74.728	71.735	93.078	77.118	75.382	72.308	92.825	77.167	75.868	72.419
+ASAM 1.0	93.098	77.983	75.47	71.554	92.969	77.947	75.154	71.268	93.143	78.058	75.724	71.695	92.797	77.805	76.222	72.034
Low Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
RL	85.745	98.603	71.162	86.714	85.451	98.463	70.768	86.192	86.472	98.74	71.522	87.63	86.865	98.95	72.36	88.594
+ASAM 1.0	88.517	99.408	73.728	91.069	88.218	99.377	73.32	90.425	88.985	99.457	73.758	91.516	88.963	99.507	74.072	91.74
SalUn	91.991	99.778	74.612	95.008	91.743	99.77	74.488	94.652	91.696	99.818	75.074	95.116	91.412	99.85	75.514	95.218
+ASAM 1.0	92.095	99.85	75.224	95.628	91.882	99.818	74.992	95.224	91.967	99.857	75.676	95.924	91.579	99.873	75.964	95.791
CIFAR100	1	4 =	SGD		1	A = ASAM 0.1 $A = A$				4 = AS	A = ASAM 1.0			A = SA	M 0 1	
High Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
L1-Sparse	74.76	5.267	61.49	63.448	75.426	5.067	60.89	63.699	73.969	6.167	60.17	61.252	77.429	7.133	62.56	65.258
+ASAM 1.0	77.86	5.733	62.99	66.903	77.648	5.7	62.29	66.213	77.126	6.367	62.02	65.117		6.2	60.83	63.051
SCRUB +ASAM 1.0	99.867 99.962	44.567 53.533	74.52 76.06	58.418 50.313	99.793 99.955	35.267 42.633	73.85 74.72	67.163 60.515	99.902	45.233 55.3	74.59 76.14	57.816 48.569	99.971	60.7 85.567	76.47 77.23	43.246 18.137
RL +ASAM 1.0	82.681 84.012	9.233 9.7	62.95 63.88	68.464 69.952	79.229	8.367 8.4	60.7 61.41	64.518 66.909	82.99 86.195	10.933 12	61.92 63.53	66.689 69.73	81.069 89.324	10.833 13.7	60.82	64.391 72.884
SalUn +ASAM 1.0	89.624 94.557	16.567 20.9	64.88 68.96	69.926 73.268	86.298 92.326	15.467 18.3	62.71 65.94	66.541 71.426	91.207	20.7 25.033	64.33 66.46	67.355 67.715	90.593	18.533 24.367	65.65 68.89	69.671 70.933
Mid Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
L1-Sparse	67.864	36.8	57.97	68.686	71.305	38.633	59.98	72.775	68.264	37.933	57.67	68.197	71.495	39.967	59.73	71.941
+ASAM 1.0	74.148	41.5	61.96	75.554	75.836	42.7	62.7	77.119	74.267	43.967	61.59	73.754	73.857	40.667	60.52	74.556
SCRUB +ASAM 1.0	99.864 99.974	81.4 85.133	74.29 75.51	76.125 73.353	99.876	76.9 77.367	72.37 74.24	79.09 80.204	99.91	83.867 85.433	73.59 75.56	73.176 73.09	99.974	90.167 97.667	75.78 77.13	68.433 61.618
RL	79.262	37.067	62.53	84.395	75.757	31.733	58.31	80.215	81.955	36.433	61.21	86.411	81.905	38.033	61.48	85.481
+ASAM 1.0	81.688	38.7	63.54	86.779	81.686	37.333	62.3	86.557	85.674	38.7	63.65	91.124	84.914	40.167	63.08	88.633
SalUn	82.383	40.733	60.46	83.056	82.4	40.9	60.9	83.377	89.581	45.333	63.46	89.768	90.205	46.867	64.8	90.495
+ASAM 1.0	91.579	48.167	66.23	92.225	88.71	45.833	64.15	89.182	94.217	50.5	66.77	93.401		52.333	67.91	92.914
Low Mem	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
L1-Sparse	62.41	91.367	55.39	53.991	68.667	96	60.25	60.34	68.421	94.2	60.67	61.47	67.229	94.967	59.33	59.014
+ASAM 1.0	66.95	94.5	59.24	58.967	73.457	96.4	63.4	66.697	70.207	96.1	61.46	62.517	72.355	96.2	62.46	65.117
SCRUB +ASAM 1.0	17.81 99.683	32.6 99.9	18.33 73.61	12.708 97.631	15.698 99.869	28.367 99.833	15.87 73.24	10.823 97.508	66.324	90.167 99.8	56.04 73.7	58.483 97.776	23.038	43.7 99.8	23.95 73.77	17.368 97.933
RL	76.376	89.233	61.34	72.4	73.283	86.5	59.57	69.711	73.495	84.2	57.63	69.677	76.79	91.733	60.62	70.55
+ASAM 1.0	78.286	90.533	62.59	74.409	73.881	87.3	59.08	69.375	82.695	89.333	63.53	80.321	83.483	94.167	64.12	78.066
SalUn	78.867	92.667	60.5	71.73	77.748	88.833	59.01	71.95	83.921	91.2	62.39	79.095	82.221	93.133	61.44	75.281
+ASAM 1.0	91.205	95.5	68.28	88.175	90.043	93.367	65.47	86.13	93.812	95.8	67.11	89.289	91.848	95.933	66.24	86.477

Table 11: Differed settings of pretrained models and their test accuracies using different A, as well as performance of retrained models w.r.t  $\mathcal{F}_{rand}$  for CIFAR-10 and Tiny-ImageNet.

Dataset, Model	lr+warmup	Batch $B$	Epoch T	W. Decay	SGD	ASAM 0.1	ASAM 1.0	SAM 0.1	Retain	Forget	Test
CIFAR10, Res18 TinyImgNt, Res34		256 256	50 200	5e-4 1e-3	93.02	93.26 62.77	93.7 62.74			92.567 59.383	

## G.3 UNLEARNING WITH STRUCTURED NOISE

We consider a noisy unlearning case where only a corrupted version of S is available, following corruptions in ImageNet-C (Hendrycks & Dietterich, 2019) to apply glass blur and snow effect to CIFAR-100 with medium severity for additional empirical verification, and report ToWs in Tab. 13: We observe that SAM continues to improve base unlearning methods with even more clear margins. This is because that structured noise applying to the images affects the dataset's signal and noise vectors ( $\varphi$  and  $\xi_i$ ), causing a corrupted dataset with worse initial signal-noise ratio, but it does not affect update dynamics and the gained results under our theoretical framework, as corrupted images are still visually recognizable, and SGD still overfits more to the added noise.

#### G.4 SAM WITH ADAM AND VIT

We also verify that our observations generalize to different base optimizers and architectures. We experiment CIFAR-100 unlearning using ViT-Small (Dosovitskiy et al., 2020) and AdamW (Loshchilov & Hutter, 2017), and summarize our priliminary results in Tab. 14. For pretraining, we use AdamW with starting Ir 0.0001, weight decay 0.05, and set patch size to 4 for

Table 12: Detailed accuracies of previous methods on Tiny-ImageNet and CIFAR-10.

TinyImageNet	1	A = 5	SGD			A = AS	AM 0.1			A = AS	AM 1.0		A = SAM 0.1			
Random $\mathcal{F}_{rand}$	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW	Retain	Forget	Test	ToW
L1-Sparse +ASAM 1.0	79.247 89.379	52.233 59.5	49.61 54.37	74.669 <b>82.753</b>	82.722 90.81	54.217 60.933	50.81 54.35	77.545 <b>82.853</b>	84.63 92.005	59.583 63.517	53.01 53.7	77.143 <b>81.168</b>	76.005 94.674	63.017 74.333	49.56 55.25	64.372 <b>75.347</b>
SCRUB +ASAM 1.0	92.112 97.965	58.117 57.717	53.65 56.94	85.793 <b>94.881</b>	94.315 98.941	60.75 61.833	54.58 58.13	86.425 <b>93.095</b>	96.268 99.521	66.5 68.333	55.01 57.66	83.457 <b>86.975</b>	99.801 99.962	88.233 97.267	58.99 61.05	<b>69.101</b> 61.704
RL +ASAM 1.0	64.504 69.356	63.233 68.733	46.59 49.22	52.668 <b>55.043</b>	67.506 73.517	66.433 72.033	47.49 50.97	53.849 <b>57.345</b>	70.309 75.88	69.883 75.617	48.16 50.38	54.424 <b>56.384</b>	75.016 81.006	73.5 79.683	49.21 50.94	56.397 <b>57.632</b>
SalUn +ASAM 1.0	69.39 75.013	68.45 74.333	50 52.65	55.735 <b>58.042</b>	70.087 77.101	68.767 75.917	49.54 53.16	55.806 <b>58.876</b>	73.207 81.039	71.783 79.233	50.12 52.89	56.721 <b>59.248</b>	82.877 88.021	81.467 87.417	53.36 54.81	<b>59.206</b> 58.998
NegGrad +ASAM 1.0	84.286 90.907	47.867 50.45	50.51 54.47	83.499 <b>91.894</b>	87.031 93.681	48.467 51.35	51.45 53.66	86.662 <b>93.094</b>	86.575 96.343	52.2 54.167	51.28 54.31	83.148 <b>93.902</b>	99.979 98.031	99.167 62.767	62.51 55.21	60.706 <b>88.59</b>
MinMax +ASAM 1.0	81.8 87.654	52.833 43.183	51.14 53.4	77.977 <b>93.426</b>	82.115 88.273	54.017 43.083	50.91 52.86	77.209 <b>93.613</b>	81.418 91.947	55.433 43.6	50.32 53.37	75.025 <b>97.617</b>	68.67 94.517	54.217 48.5	46.99 53.72	61.615 <b>96.466</b>
CIFAR10								A = ASAM 1.0								
CHIMIN		A = S	SGD			A = AS	AM 0.1		1	A = AS	AM 1.0			A = SA	M 0.1	
Random $\mathcal{F}_{rand}$	Retain	A = S	Test	ToW	Retain	A =AS. Forget	AM 0.1 Test	ToW	Retain	A =AS. Forget	AM 1.0 Test	ToW	Retain	A =SA Forget	M 0.1 Test	ToW
	Retain   86.467   91.438			ToW 85.12 <b>90.352</b>	Retain 89.06 91.674			ToW 87.688 <b>91.087</b>	Retain 86.683 90.938			ToW 84.811 <b>89.268</b>	Retain 89.462 90.886			ToW 87.144 <b>88.792</b>
Random $\mathcal{F}_{rand}$ L1-Sparse	86.467	Forget 82.967	Test 82.25	85.12	89.06	Forget 85.567	Test 84.45	87.688	86.683	Forget 83.467	Test 82.11	84.811	89.462	Forget 87.133	Test 84.82	87.144
Random F <sub>rand</sub> L1-Sparse +ASAM 1.0 SCRUB	86.467 91.438 90.767	Forget 82.967 88.333 86.033	Test 82.25 87.23 86.27	85.12 <b>90.352</b> 90.739	89.06 91.674 68.205	Forget 85.567 87.767 67.367	Test 84.45 87.24 66.75	87.688 <b>91.087</b> 63.466	86.683 90.938 80.193	Forget 83.467 88.7 78.933	Test 82.11 86.94 77.97	84.811 <b>89.268</b> 77.95	89.462 90.886	Forget 87.133 88.633	Test 84.82 86.43	87.144 <b>88.792</b> 6.089
Random F <sub>rand</sub> L1-Sparse +ASAM 1.0 SCRUB +ASAM 1.0 RL	86.467 91.438 90.767 99.6 92.774	Forget 82.967 88.333 86.033 95.167 86.6	Test 82.25 87.23 86.27 92.65 87.22	85.12 90.352 90.739 97.2 93.186	89.06 91.674 68.205 99.621 90.569	Forget 85.567 87.767 67.367 96.5 84.2	Test 84.45 87.24 66.75 93.15 85.17	87.688 <b>91.087</b> 63.466 <b>96.39</b> 91.02	86.683 90.938 80.193 99.807 91.445	Forget 83.467 88.7 78.933 98.2 84.133	Test 82.11 86.94 77.97 93.38 85.81	84.811 <b>89.268</b> 77.95 <b>95.078</b> 92.591	89.462 90.886 15.11 99.631 88.736	Forget 87.133 88.633 14.2 98.467 82.533	Test 84.82 86.43 15 93.16 84.12	87.144 <b>88.792</b> 6.089 <b>94.435</b> 89.524
Random F <sub>rand</sub> L1-Sparse +ASAM 1.0 SCRUB +ASAM 1.0 RL +ASAM 1.0 SalUn	86.467   91.438   90.767   99.6   92.774   93.295   96.94	Forget 82.967 88.333 86.033 95.167 86.6 87.733	Test 82.25 87.23 86.27 92.65 87.22 87.66	85.12 90.352 90.739 97.2 93.186 93.138 98.095	89.06 91.674 68.205 99.621 90.569 93.262 95.726	Forget 85.567 87.767 67.367 96.5 84.2 87.233 87.6	Test 84.45 87.24 66.75 93.15 85.17 88.31 89.02	87.688 91.087 63.466 96.39 91.02 94.187 97.052	86.683   90.938   80.193   99.807   91.445   95.098	Forget 83.467 88.7 78.933 98.2 84.133 89.033 88.733	Test 82.11 86.94 77.97 93.38 85.81 89.44 89.35	84.811 <b>89.268</b> 77.95 <b>95.078</b> 92.591 <b>95.512</b> 96.598	89.462   90.886   15.11   99.631   88.736   92.588   96.612	Forget 87.133 88.633 14.2 98.467 82.533 86.567 89.867	Test 84.82 86.43 15 93.16 84.12 87.4	87.144 88.792 6.089 94.435 89.524 93.206

Table 13: Unlearning with ImageNet-C corruptions on CIFAR-100.

Glass Blur		A=5	SGD			A=A	SAM	
Method	High	Mid	Low	AVG	High	Mid	Low	AVG
NG	67.760	78.824	75.931	74.172	76.152	85.534	82.556	81.414
+ASAM	73.565	80.253	84.086	79.301	74.993	86.567	86.296	82.619
SharpMinMax	66.110	76.852	73.387	72.116	66.837	79.023	78.435	74.765
+ASAM	75.327	89.859	79.104	81.430	74.089	92.737	84.921	83.916
Snow		A=5	SGD			A=A	SAM	
Snow Method	High	A=S Mid	SGD Low	AVG	High	A=A Mid	SAM Low	AVG
	   High   77.394			AVG 81.306	High 75.041			AVG 82.768
Method		Mid	Low			Mid	Low	
Method NG	77.394	Mid 83.328	Low 83.196	81.306	75.041	Mid 86.424	Low 86.838	82.768

ViT-Small on CIFAR-100. Other experiment settings are unchanged. For unlearning, we have unlearn lr 0.0006 for NegGrad and for Sharp MinMax. Adam demands much smaller lr than SGD and is more sensitive to unlearn lr tuning. ViTs perform worse than ResNets on smaller datasets (test accuracies of pretrained models are 57%).

#### H COMPLETE VISUALIZATIONS

In this section, we provide complete visualizations of feature space and loss landscapes of pretrained models, NegGrad unlearned models, and Sharp MinMax unlearned models, comparing SGD with SAM across all memorization levels. The observations are generally consistent across memorization levels, with  $\mathcal{F}_{high}$  being more noticeable.

## H.1 LOSS LANDSCAPE

Inspired by Wu et al. (2017), we quantify the flatness by basin ratio, which is the percentage of perturbed losses whose deviation from original loss  $\leq 0.5 \cdot \text{stddev}$ . Fig. 5 shows loss landscapes of SAM and SGD before and after unlearning on  $\mathcal{D}_{\text{test}}$  and  $\mathcal{F}_{\text{high}}$ . We observe SAM has higher basin ratios (flatter landscape) than SGD for pretrained model and MinMax unlearned model as expected. Surprisingly, SGD can become flatter after unlearning. We conjecture that the gradient ascent might

Table 14: Unlearning with ViT-Small and AdamW on CIFAR-100.

		A=	SGD		A=ASAM						
	High	Mid	Low	AVG	High	Mid	Low	AVG			
NG	80.445	82.854	84.385	82.561	78.750	82.223	86.767	82.580			
+ASAM	82.880	83.084	83.402	83.122	82.839	81.354	87.507	83.900			
SharpMinMax	14.794	42.055	95.222	50.690	14.279	42.017	94.833	50.376			
+ASAM	76.343	95.573	103.372	91.763	76.664	93.966	105.868	92.166			

be implicitly regularizing SGD which had more overfitting than SAM during pretraining. We leave the further characterization of loss landscapes to future work.

## H.2 FEATURE VISUALIZATION

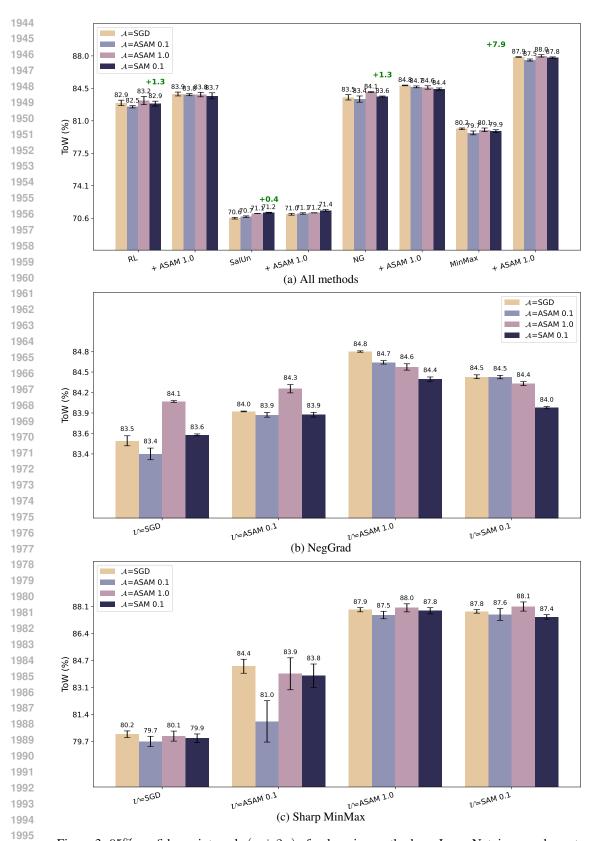


Figure 3: 95% confidence intervals ( $\mu \pm 2\sigma$ ) of unlearning methods on ImageNet, in accordance to Tab. 1 and Tab. 3. We run each setting three times with different seeds and compute the statistical significance. SAM consistently improves base  $\mathcal{U}$ , and we observe ASAM 1.0 to bring largest improvement steadily.

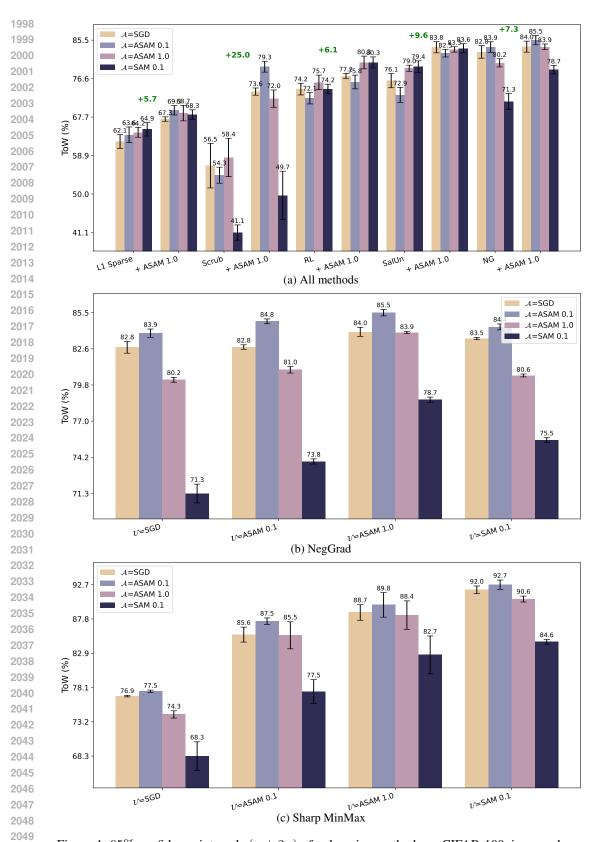


Figure 4: 95% confidence intervals ( $\mu \pm 2\sigma$ ) of unlearning methods on CIFAR-100, in accordance to Tab. 1 and Tab. 3. We run each setting three times with different seeds and compute the statistical significance. SAM not only improves ToW of the based methods, but also more robust against variance than SGD.

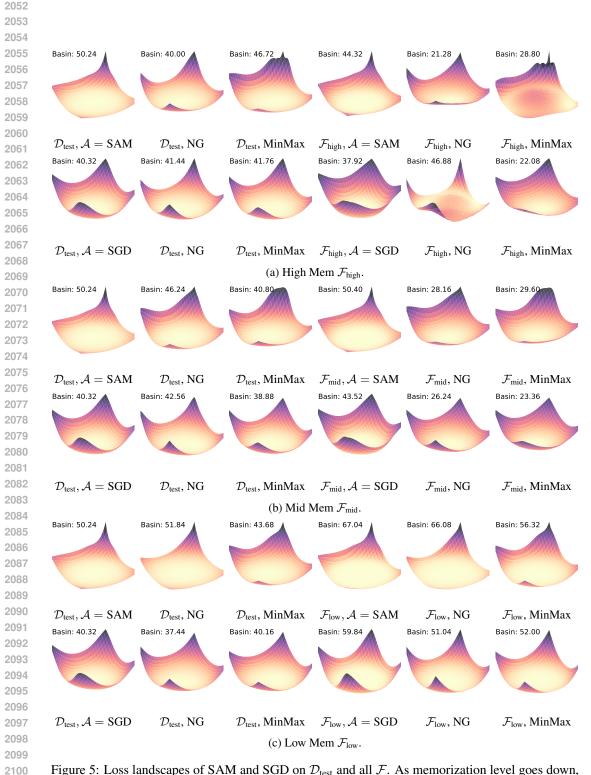


Figure 5: Loss landscapes of SAM and SGD on  $\mathcal{D}_{test}$  and all  $\mathcal{F}$ . As memorization level goes down,  $\mathcal{F}$  becomes easier to unlearn and SGD shows less to no "regularizing" effect as we have discussed on  $\mathcal{F}_{high}$ . The general trend preserves with decreasing memorization levels and SAM is generally flatter before and after unlearning.

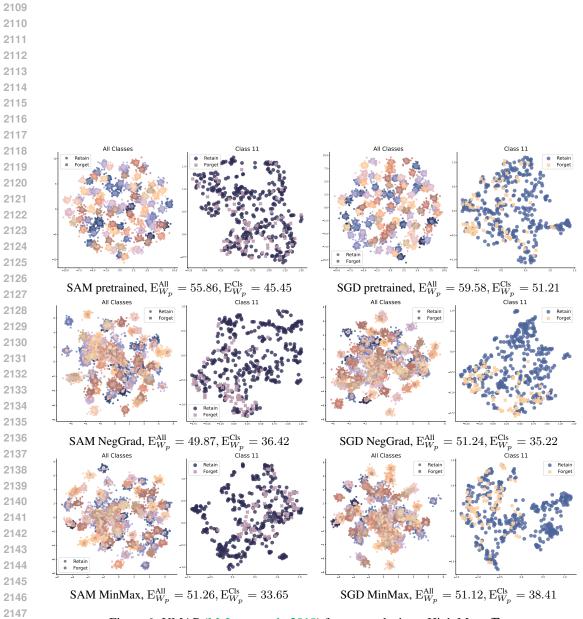


Figure 6: UMAP (McInnes et al., 2018) feature analysis on High Mem  $\mathcal{F}_{high}$ .

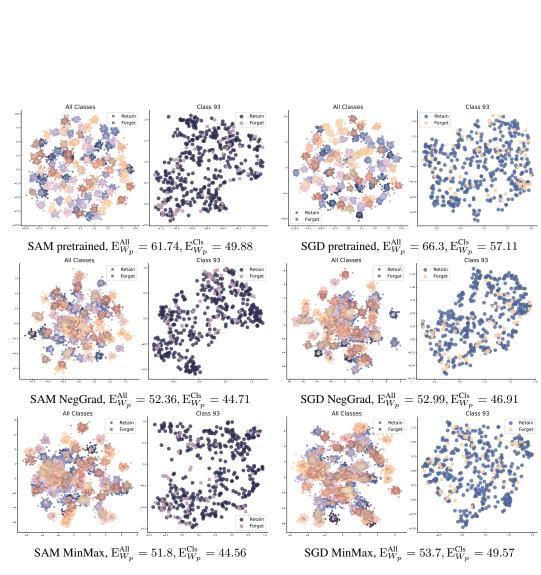


Figure 7: UMAP (McInnes et al., 2018) feature analysis on Mid Mem  $\mathcal{F}_{mid}$ .

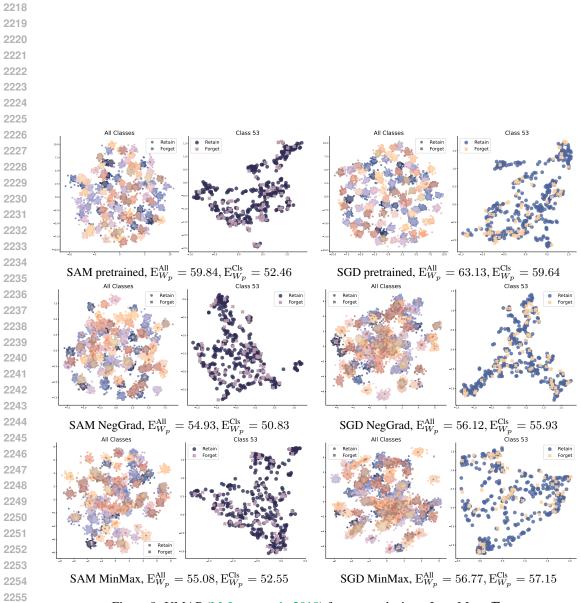


Figure 8: UMAP (McInnes et al., 2018) feature analysis on Low Mem  $\mathcal{F}_{low}$ .