

INDUCING UNCERTAINTY ON OPEN-WEIGHT MODELS FOR TEST-TIME PRIVACY IN IMAGE RECOGNITION

Anonymous authors

Paper under double-blind review

ABSTRACT

A key concern for AI safety remains understudied in the machine learning (ML) literature: how can we ensure users of ML models do not leverage predictions on incorrect personal data to harm others? This is particularly pertinent given the rise of open-weight models, where simply masking model outputs does not suffice to prevent adversaries from recovering harmful predictions. To address this threat, which we call *test-time privacy*, we induce maximal uncertainty on protected instances while preserving accuracy on all other instances. Our proposed algorithm uses a Pareto optimal objective that explicitly balances test-time privacy against utility. We also provide a certifiable approximation algorithm which achieves (ϵ, δ) guarantees without convexity assumptions. We then prove a tight bound that characterizes the privacy-utility tradeoff that our algorithms incur. Empirically, our method obtains at least $> 3\times$ stronger uncertainty than pretraining with marginal drops in accuracy on various image recognition benchmarks. Altogether, this framework provides a tool to guarantee additional protection to end users.

1 INTRODUCTION

Data privacy is increasingly important for large-scale machine learning (ML), where models are often trained on sensitive user instances (European Parliament, 2016). Furthermore, open-weight image recognition models, where users have access to the model parameters and architecture, have proliferated (TorchVision, 2016; Google, 2023; Microsoft, 2024).

Yet, there has been little work done to address privacy threats to ML models due to incorrect personal data, especially data which are public such as images posted to public forums. Concretely, suppose a model provider trains an open-weight medical imaging model f which classifies skin images as harmless ailments like “Benign Keratosis” or serious diseases like “Melanoma” (Sun et al., 2016). Next, a health insurance company scrapes images from public forums to build risk profiles. Then, this health insurance company downloads the open-weight model f to automatically screen images for potential health liabilities. In particular, a person p posts a photo of a harmless birthmark to a public health forum to ask a question. During the upload, a compression error causes the image file to become corrupted, severely distorting the birthmark. This results in an image x_p . When the health insurance company feeds x_p , scraped from the public forum, into f , it confidently classifies x_p as “Melanoma”. This erroneous classification is then automatically added to person p ’s risk profile, resulting in person p being unfairly denied coverage.¹ We call this threat model *test-time privacy* (TTP), and make this concrete in Fig. 1.² **This privacy threat model is inspired by definitions of privacy which correspond to protecting a user from unfair interference or intrusion (Merriam-Webster, 2022). This differs from settings in privacy which mainly protect sensitive information.**

We discuss why existing solutions, like unlearning (Sekhari et al., 2021) or differential privacy (Dwork et al., 2006) do not suffice to solve this problem in Sec. 2. Furthermore, naive solutions like masking model outputs do not work for open-weight image recognition models, since the model parameters

¹Recently, there has been significant progress in building effective image recognition models for skin disease, making this problem pertinent (Yang et al., 2018); (Liu et al., 2025).

²We provide additional test-time privacy examples beyond medical classifiers in Appx. E.

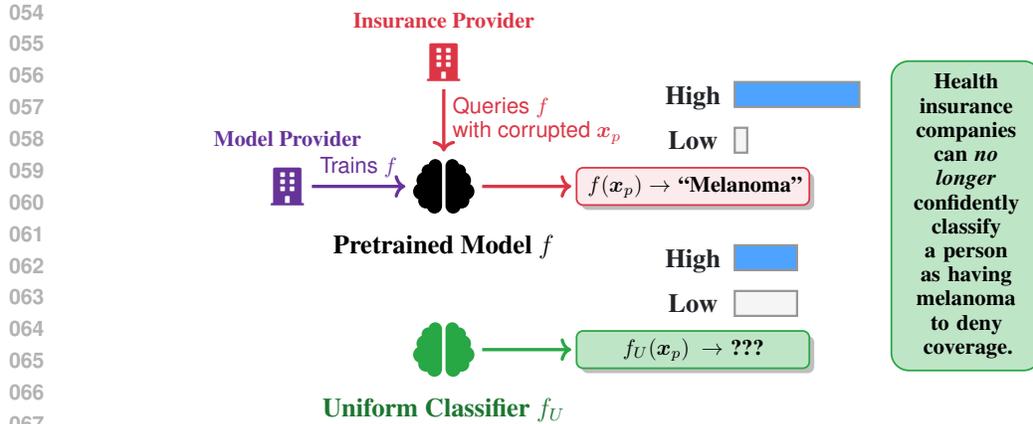


Figure 1: An adversary, like a health insurance company (🏢), can query a pretrained model f (🧠) and use its outputs to make harmful decisions. However, after running our algorithm, the new model f_U (🌱) provides maximal uncertainty, protecting against such an adversary.

and architecture are available to the model user. An adversary could simply remove such a mask. To make this clear, we comprehensively detail our threat model as a security game in Appx. A and provide various motivating attacks in Appx. B. Therefore, we ask the following research question:

Can we ensure test-time privacy against adversaries with access to an open-weight model?

To do so, we argue it suffices to have uniform model outputs over the protected instances. That way, a data controller can only guess at the prediction. Thus, we revisit inducing maximal uncertainty over a dataset (Pereyra et al., 2017). Furthermore, we want to obtain high performance on all other instances as well. In particular, we answer our research question affirmatively, providing:

- A method to finetune a pretrained model with a Pareto optimal objective, rendering the model maximally uncertain over protected instances while preserving accuracy on others.
- Several principled (ϵ, δ) -certified full-batch and sequential algorithms which approximate the Pareto objective, derived without assumptions of convexity.
- A theoretical analysis of the privacy-utility tradeoff that our algorithms incur, establishing a tight, non-vacuous bound.
- Empirical studies on image recognition models like ResNet50 (He et al., 2016) trained on datasets like CIFAR100 (Krizhevsky et al., 2009), observing that our algorithms maintain high uniformity on protected instances while guaranteeing excellent utility on the rest.

Following the literature on privacy, we focus on protecting a subset of the training data. However, as detailed in Sec. 4 and Appx. A, our setup and algorithms can also work for corrupted test instances. The code for our experiments is available for reproducibility at <https://tinyurl.com/testtimeprivacy>.

2 RELATED WORK

Data Privacy: Data leakage is a persistent danger for large information systems (Al-Rubaie and Chang, 2019). In the context of ML, data privacy is ubiquitous (Fredrikson et al., 2015); (Song et al., 2017); (Yeom et al., 2018). Approaches to privacy-preserving ML include differential privacy (DP) (Dwork et al., 2006); (Balle and Wang, 2018); (Amin et al., 2024), homomorphic encryption (Brakerski et al., 2014); (Lee et al., 2022); (Aono et al., 2017), and model obfuscation (Zhou et al., 2023). Notably, these methods protect against various privacy violations, like reconstruction attacks (Dinur and Nissim, 2003) due to failures of anonymization (Li et al., 2012). However, existing methods do not prevent confidently correct classification, and thus fail to protect against the attacks we consider in our setting. For example, if x_p is a corrupted medical image record, an adversary may not be able to use a DP model f to recover the record x_p exactly, but they can still produce a confident prediction of e.g. "Melanoma" to harm person p .

108 A dominant viewpoint in the privacy community is that a model f working as expected does not
 109 constitute a privacy violation, e.g. correctly predicting “Melanoma” for the corrupted medical
 110 image x_p , as it has learned something underlying about nature (Mcsherry, 2016); (Bun et al., 2021).
 111 Furthermore, the privacy leakage occurred when x_p became public (Kamath, 2020). This view misses
 112 the point: ML models are often trained and applied on freely available data. For example, training
 113 data could be scraped from the web or social media platforms. Subsets of this data can be obsolete,
 114 corrupted, or confidential. With such data as input, model f presents a clear and present danger
 115 for AI safety which differential privacy falls short of addressing, as the above example showed. In
 116 parallel, as humans, we learn to not act upon certain kinds of knowledge. For example, when we read
 117 confidential documents or learn that previously obtained knowledge is incorrect, we are not allowed
 118 to share or act on this knowledge.

119 **Unlearning:** A related subfield is machine unlearning, which is inspired by the right to be forgotten
 120 (RTBF), mandating that ML model providers delete user data upon request (European Parliament,
 121 2016). In practice, model providers must remove user data and its effects from trained models and
 122 algorithms. Unlearning methods usually do so by approximating (and evaluating performance against)
 123 the model retrained from scratch without the protected user data (Sekhari et al., 2021); (Bourtole
 124 et al., 2021); (Kurmanji et al., 2023).

125 However, while unlearning helps model providers comply with the RTBF, it cannot protect against
 126 attacks within our threat model. Specifically, recent unlearning research has established that data
 127 in the support of the training distribution will likely still be confidently predicted with the same
 128 prediction as before, even after using state-of-the-art algorithms (or even after applying exact retrain-
 129 from-scratch algorithms) to unlearn them (Zhao et al., 2024). That is, denoting the pretrained model
 130 as f and the unlearned model as f_u , for typical training instances, it holds that $f = f_u$.

131 To make clear why unlearning does not solve our problem, recall the example of a model f trained
 132 on skin images to predict disease. This time, to remove the corrupted medical image x_p from f ,
 133 person p invokes the RTBF. Thus, the data controller for f unlearns x_p , yielding f_u . But, even after
 134 unlearning, any medical insurance company can still access the publicly available x_p and obtain
 135 $f_u(x_p)$. But, $f_u(x_p) = f(x_p)$, and thus the medical insurance company *incorrectly* labels person
 136 p as high risk for medical coverage. Thus, the unfair and dangerous scenario for person p remains.
 137 This holds similarly for unlearning methods which deal with corrupted or obsolete data, as they still
 138 do not aim to reduce confidence in the final prediction (Schoepf et al., 2025).

139 *Differences from Unlearning:* Importantly, what we propose is **not an unlearning algorithm**,
 140 which would need to be aligned with the goals of unlearning (and indistinguishability from retrain-
 141 from-scratch). Instead, we aim to address an entirely new threat scenario—test time privacy—that
 142 unlearning cannot solve, which we detail in Appx. A. For example, indistinguishability from retrain
 143 is inconsequential in our threat model. Furthermore, we also consider corrupted test examples, unlike
 144 unlearning which focuses only on the training dataset. Finally, what may constitute a privacy violation
 145 in unlearning, e.g. revealing that an instance is in the forget set via a membership inference attack
 146 (Shokri et al., 2017) does *not* constitute a violation in our threat model. This holds similarly for other
 147 threat models; for example, reconstruction attacks which lead to recovery of x_p (Dinur and Nissim,
 148 2003) or adversarial attacks which lead to misclassification of x_p (Goodfellow et al., 2015) are not
 149 violations in our threat model, as explained in Appx. A.

150 Still, the privacy guarantees that we provide in the threat model of test-time privacy are complementary
 151 to the guarantees that unlearning can provide. However, related work has focused heavily on
 152 unlearning—we fill this gap by presenting a framework for test-time privacy.

153 **Additional Related Works:** Due to space constraints, we provide an additional related works section
 154 in Appx. C. Critically, we describe how differentially privacy methods like private aggregation of
 155 teacher ensembles (PATE) (Papernot et al., 2018) or label differential privacy (LabelDP) (Ghazi et al.,
 156 2021) differ from our setting, how label model inversion attacks (Zhu et al., 2019) relate to our threat
 157 model, and also why misclassification-based methods for unlearning (Cha et al., 2024) are suboptimal
 158 for addressing our threat model.

3 TEST-TIME PRIVACY THREAT MODEL

Notation: Let $\mathcal{X} \subset \mathbb{R}^d$ be a sample space and let $\mathcal{Y} \subset \mathbb{R}^o$ be a label space. Denote $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ as the space of feature-label pairs. Let \mathcal{Z}^n be the n -fold Cartesian product of \mathcal{Z} such that a dataset $\mathcal{D} \subset \mathcal{Z}^n$ is a collection of n feature-label pairs. Then, the i th instance is denoted as $\mathcal{D}^{(i)}$ with its feature in \mathcal{X} being $\mathcal{D}^{(i,\mathcal{X})}$ and label being $\mathcal{D}^{(i,\mathcal{Y})}$. Following the unlearning literature, we subset \mathcal{D} as a “forget set” \mathcal{D}_f , containing instances to protect, and a retain set $\mathcal{D}_r = \mathcal{D} \setminus \mathcal{D}_f$. Then, suppose we have a (randomized) learning algorithm $\mathcal{A} : \mathcal{Z}^n \rightarrow \mathcal{W}$, where $\mathcal{W} \subset \mathbb{R}^z$ is a parameter space. Let the set of hypotheses parameterized with respect to this parameter space be $\mathcal{H}_{\mathcal{W}}$. Let $f_w \in \mathcal{H}_{\mathcal{W}}$ be the hypothesis parameterized by $w \in \mathcal{W}$, defined as $f_w : \mathcal{X} \rightarrow \Delta_{|\mathcal{Y}|}$, where $\Delta_{|\mathcal{Y}|}$ is the probability simplex $\{p_1, \dots, p_\gamma : p_i \geq 0, \sum_{i=1}^{|\mathcal{Y}|} p_i = 1\}$. When \mathbf{A} is a matrix, $\|\mathbf{A}\|_2$ is the 2 operator norm. When \mathbf{v} is a vector, $\|\mathbf{v}\|_2$ is the ℓ_2 norm. Furthermore, let $\lambda_{\min}(\mathbf{A})$ denote the minimum eigenvalue of \mathbf{A} . If we have an objective $f_A + f_B$, we denote its gradient evaluated at w as $\nabla_{w,A,B}$ and Hessian as $\mathbf{H}_{w,A,B}$. Finally, when for sets $\mathcal{S}, \mathcal{F} \subset \mathcal{N}$ and \mathcal{R} and mechanisms $\mathcal{M}, \mathcal{M}' : \mathcal{N} \rightarrow \mathcal{R}$, we have $\Pr(\mathcal{M}(\mathcal{S}) \in \mathcal{R}) \leq e^\epsilon \Pr(\mathcal{M}'(\mathcal{F}) \in \mathcal{R}) + \delta$ and $\Pr(\mathcal{M}'(\mathcal{S}) \in \mathcal{R}) \leq e^\epsilon \Pr(\mathcal{M}(\mathcal{F}) \in \mathcal{R}) + \delta$, we will denote $\mathcal{M}(\mathcal{S}) \approx_{\epsilon,\delta,\mathcal{R}} \mathcal{M}'(\mathcal{F})$. We provide a symbol table in Appx. N.

Following recent work (Baig and Pietrzak, 2025), we define our threat model as an informal *test-time privacy (TTP) game*. This section summarizes the full model detailed in Appx. A.

Firstly, the game involves three actors: a **data corrupter** (c) who creates a “forget set” \mathcal{D}_f of corrupted instances; a benign **model provider** (τ) who releases a model; and a **TTP adversary** (ν) who aims to use the model to get a confident, harmful prediction m on instances $x \in \mathcal{D}_f$. Next, the provider τ trains a model f , is notified of the corrupted set \mathcal{D}_f , and then runs an algorithm \mathcal{G} to release a protected model \hat{f} . We make two key assumptions:

Assumption 3.1 (Open-Weight Access). *The adversary ν has full access to \hat{f} ’s architecture and weights.*

Assumption 3.2 (τ -Limited Knowledge). *The provider τ knows \mathcal{D}_f but does not know the specific harmful label m the adversary seeks, nor do they know the adversary ν .*

The adversary ν wins if they can leverage open-weight access to \hat{f} to *confidently* recover the sensitive prediction m . The provider τ wins if \hat{f} leaves ν uncertain about m . Given this, together with Asm. 3.2, the provider’s optimal defense is to have \hat{f} output maximal uncertainty (i.e. a uniform distribution) on \mathcal{D}_f while maintaining accuracy on all other data. Critically, Asm. 3.1 renders simple defenses like output masking useless, as the adversary ν could just remove the mask. Therefore, a successful defense must perturb the model weights in a non-invertible manner, which motivates our approach in Sec. 4.

Finally, our formulation is general. While we may define \mathcal{D}_f in terms of training data, this is done without loss of generality; the forget set can consist of any corrupted instances, whether from the training or test set. We provide simple TTP attacks to further motivate this threat model in Appx. B.

4 APPROACHES AND ALGORITHMS

In order to prevent test-time privacy violations, it suffices to have the model output a uniform distribution over the forget set, rendering the model maximally uncertain. Then, an adversary can only guess at the original sensitive prediction.³ We also would like to preserve retain set accuracy; to that end, we present an algorithm that finetunes the pretrained model with a Pareto optimal objective. To make this algorithm concrete, we define a *uniform learner*, which we prove to exist in many common hypothesis classes. Then, we use this concept in order to construct a Pareto objective.

4.1 THE FINETUNING PARETO LEARNER

For a dataset $\mathcal{D} \subset \mathcal{Z}^n$, we denote the pretrained model as $\mathcal{A}(\mathcal{D})$. Then, to make $\mathcal{A}(\mathcal{D})$ uniform over the forget set, we introduce our core concept of a *uniform learner*:

³Please see Appx. A for a clear characterization of why this is optimal within our threat model, and why we can consider $\mathcal{D}_f \subset \mathcal{D}$, where \mathcal{D} is a training dataset, without loss of generality.

Definition 4.1 (Uniform learner). *Suppose we have a (randomized) learning algorithm $\mathcal{K} : \mathcal{Z}^n \rightarrow \mathcal{W}$ that, given $\mathcal{D} \subset \mathcal{Z}^n$, yields the parameter $\mathcal{K}(\mathcal{D}) = \mathbf{w}_{\mathcal{D}}$. We say \mathcal{K} is a uniform learner if $\forall \mathcal{D} \in \mathcal{Z}^n$, $\mathbf{w}_{\mathcal{D}}$ parametrizes $f_{\mathbf{w}_{\mathcal{D}}} \in \mathcal{H}_{\mathcal{Y}}$ and satisfies:*

$$\|f_{\mathbf{w}_{\mathcal{D}}} - \underbrace{\left(\frac{1}{|\mathcal{Y}|}, \dots, \frac{1}{|\mathcal{Y}|}\right)}_{|\mathcal{Y}| \text{ times}}\|_{\infty, \mathcal{X}} = 0. \quad (1)$$

That is, \mathcal{K} is a uniform learner if its parameterized outputs yield the uniform distribution $U[0, |\mathcal{Y}|]$ for all inputs across all datasets. We define this as a learning algorithm for full generality to handle e.g. neural networks with nonlinear transformations in their last layer. Furthermore, \mathcal{K} exists in many common hypothesis spaces, proved in Appx. H.2:

Proposition 4.2. *Suppose we have a hypothesis space $\mathcal{H}_{\mathcal{Y}}$ consisting of functions where the ultimate layer is an affine transformation and the outputs are passed through a softmax. Let \mathcal{K} be a uniform learner. Then, $f_{\mathcal{K}(\mathcal{D})} \in \mathcal{H}_{\mathcal{Y}} \forall \mathcal{D} \subset \mathcal{Z}^n$.*

Proof Sketch: Setting the weights in the ultimate affine layer to 0 yields uniform outputs.

Most classifiers built and deployed in ML recently, including multilayer perceptrons (MLP), residual networks (ResNets) (He et al., 2016), and transformers (Vaswani et al., 2017) satisfy the premise of Prop. 4.2, making it widely applicable.

Next, we assume \mathcal{A} and \mathcal{K} are obtained through empirical risk minimization (ERM) to a local or global minima. That is, let $\mathcal{A}(\mathcal{D}) = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D})$ and $\mathcal{K}(\mathcal{D}) = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D})$, where $\mathcal{L}_{\mathcal{A}}$ penalizes incorrect classification and $\mathcal{L}_{\mathcal{K}}$ penalizes a lack of uniformity in model outputs. One choice of $\mathcal{L}_{\mathcal{K}}$ is the KL divergence (Kullback and Leibler, 1951) between the softmax outputs and the uniform distribution. This loss has been previously used to penalize highly confident classifier predictions (Pereyra et al., 2017); we thus adapt this loss to our setting. Furthermore, by Prop. 4.2, this loss can be completely minimized over \mathcal{W} .

Critically, we seek the optimal tradeoff between uniformity over the forget set and utility over the retain set. That is, we should produce a learner that is Pareto optimal with respect to $\mathcal{L}_{\mathcal{K}}$ and $\mathcal{L}_{\mathcal{A}}$. This learner can be characterized as follows:

Proposition 4.3. *Let $\theta \in (0, 1)$. Fix $\mathcal{D} \subset \mathcal{Z}^n$ and consider the forget set $\mathcal{D}_f \subset \mathcal{D}$ and the retain set $\mathcal{D}_r = \mathcal{D} \setminus \mathcal{D}_f$. Then, if $\mathcal{M}_{\theta}(\mathcal{D}) = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \theta \mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r)$ is a global minimizer, it is globally Pareto optimal with respect to $\mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f)$ and $\mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r)$. Similarly, if $\mathcal{M}_{\theta}(\mathcal{D})$ is a local minimizer, it is locally Pareto optimal.*

Proof Sketch: This holds by contradiction. If the solution to the minimized objective was not globally (locally) Pareto optimal, since $\theta \in (0, 1)$, it could not be the global (local) minimizer. See Appx. H.3 for a full proof. Definitions of Pareto optimality are included in Appx. H.3 as well.

As shown in Prop. 4.3, \mathcal{M}_{θ} yields a parameter that, given $\theta \in (0, 1)$, presents a Pareto optimal tradeoff between uniformity over the forget set and utility over the retain set. One can adjust θ to vary over many Pareto optimal solutions, yielding different tradeoffs between uniformity and utility. This yields Alg. 1, in which we finetune a pretrained model by using it as initialization for $\mathcal{M}_{\theta}(\mathcal{D})$.

4.2 THE CERTIFIED PARETO LEARNER

While the aforementioned algorithm in 4.1 provably guarantees an optimal tradeoff, so long as its objective is minimized, we would also like to make it *certified*, obtaining a certificate that a third party can inspect to verify test-time privacy. Thus, to design a certified approximation algorithm, we take inspiration from certified unlearning (Zhang et al., 2024), which aims to add a small amount of structured noise such that the pretrained model becomes indistinguishable from the retrained model. In our setting, we would like to make the pretrained model indistinguishable from the solution to the Pareto objective. To do so, we define a new notion of (ε, δ) -indistinguishability and use this definition to design a certifiable algorithm, with results in Sec. 6.

Firstly, to motivate our definition, recall the definition of differential privacy (Dwork et al., 2006):

Definition 4.4 ((ε, δ) -differential privacy). *Suppose we have privacy budgets $\varepsilon \in (0, 1)$ and $\delta > 0$. A randomized algorithm $\mathcal{M} : \mathcal{Z}^n \rightarrow \mathcal{W}$ satisfies (ε, δ) -differential privacy if $\forall \mathcal{T} \subset \mathcal{W}$ and $\forall \mathcal{D}, \mathcal{D}' \in \mathcal{Z}^n$ s.t. $\|\mathcal{D} - \mathcal{D}'\|_1 \leq 1$:*

Algorithm 1 \mathcal{M}_θ Finetuning

Require: Dataset \mathcal{D} ; forget set \mathcal{D}_f ; pretrained model $w^* = \mathcal{A}(\mathcal{D})$; uniformity-utility tradeoff coefficient θ ; e epochs.
 Use w^* as initialization for the Pareto learner $\mathcal{M}_\theta(\mathcal{D})$.
 Optimize the Pareto learner $\mathcal{M}_\theta(\mathcal{D})$ for e epochs with e.g. SGD to yield w^- .
return w^- .

$$\mathcal{M}(\mathcal{D}) \approx_{\varepsilon, \delta, \mathcal{T}} \mathcal{M}(\mathcal{D}'). \quad (2)$$

This guarantees that the algorithm \mathcal{M} applied on a dataset is statistically indistinguishable from the same algorithm applied on all datasets different by one instance. One can leverage this definition to formalize certified unlearning (Sekhari et al., 2021):

Definition 4.5 ((ε, δ) -certified unlearning). *Suppose we have privacy budgets $\varepsilon \in (0, 1)$ and $\delta > 0$. Consider $\mathcal{D} \subset \mathcal{Z}^n$ and let $\mathcal{D}_f \subset \mathcal{D}$ be the forget set to be unlearned, and $\mathcal{D}_r = \mathcal{D} \setminus \mathcal{D}_f$ be the retain set. $\mathcal{U} : \mathcal{Z}^n \times \mathcal{Z}^n \times \mathcal{W} \rightarrow \mathcal{W}$ is an (ε, δ) -certified unlearning algorithm if $\forall \mathcal{T} \subset \mathcal{W}$, we have:*

$$\mathcal{U}(\mathcal{D}, \mathcal{D}_f, \mathcal{A}(\mathcal{D})) \approx_{\varepsilon, \delta, \mathcal{T}} \mathcal{A}(\mathcal{D}_r). \quad (3)$$

This formalizes making $\mathcal{A}(\mathcal{D})$ indistinguishable from $\mathcal{A}(\mathcal{D}_r)$. In light of Def. 4.5, we seek to make $\mathcal{A}(\mathcal{D})$ indistinguishable from $\mathcal{M}_\theta(\mathcal{D})$. Thus, we provide the following new definition:

Definition 4.6 ($(\varepsilon, \delta, \theta)$ -certified Pareto learner). *Suppose we have privacy budgets $\varepsilon \in (0, 1)$ and $\delta > 0$ with $\mathcal{D} \subset \mathcal{Z}^n$. Let $\mathcal{D}_f \subset \mathcal{D}$ be the forget set and let $\mathcal{D}_r = \mathcal{D} \setminus \mathcal{D}_f$ be the retain set. Suppose we have $\theta \in (0, 1)$ and $\mathcal{M}_\theta(\mathcal{D}) = \operatorname{argmin}_{w \in \mathcal{W}} \theta \mathcal{L}_{\mathcal{K}}(w, \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(w, \mathcal{D}_r)$, where \mathcal{K} is a uniform learner and \mathcal{A} is a learning algorithm both obtained through ERM. An algorithm $\mathcal{G} : \mathcal{Z}^n \times \mathcal{Z}^n \times \mathcal{W} \rightarrow \mathcal{W}$ is a $(\varepsilon, \delta, \theta)$ -certified Pareto learner if $\forall \mathcal{T} \subset \mathcal{W}$:*

$$\mathcal{G}(\mathcal{D}, \mathcal{D}_f, \mathcal{A}(\mathcal{D})) \approx_{\varepsilon, \delta, \mathcal{T}} \mathcal{M}_\theta(\mathcal{D}). \quad (4)$$

Discussion: Qualitatively, the conditions in Def. 4.6 mean that the model obtained by algorithm \mathcal{G} is statistically indistinguishable from a model that is Pareto optimal between utility over the retain set and uniformity over the forget set. Here, we consider the classical setting of $\varepsilon \in (0, 1)$.⁴ Finally, note that satisfying Def. 4.5 and Def. 4.6 together is not possible for forget sets which overlap; thus, a model provider should adopt whichever approach corresponds to their threat model.

One way we can design an algorithm which satisfies Def. 4.6 is by taking a Newton step towards the Pareto model and applying structured Gaussian noise; this yields Alg. 2, which is certifiable as proved in Appx. G. Using local convex approximation (Nocedal and Wright, 1999), in which we add a regularization term to the objective of the Pareto learner, we design Alg. 2 without any assumptions of convexity on the component loss functions.

In addition, Alg. 2 requires inverting a Hessian, which is computationally infeasible for practical neural networks e.g. ResNets, even after employing conjugate gradient methods (Nocedal and Wright, 1999) and Hessian vector product techniques (Pearlmutter, 1994). To resolve this issue, we also propose a derived Alg. 3 in Appx. G, which computes an efficient estimator for the inverse Hessian (Agarwal et al., 2016). Furthermore, this algorithm does not assume convergence to a local minima for $\mathcal{A}(\mathcal{D})$, handling e.g. early stopping. An online version is presented in Appx. I as Alg. 4. While Alg. 2, 3 and 4 have more hyperparameters than Alg. 1, they offer a certificate which can be used to verify use of our method by a third party; we present ways to reduce hyperparameters in Appx. J.

5 THEORETICAL ANALYSIS

In what follows, we aim to analyze various properties of Alg. 1 and Alg. 2 to understand how to appropriately choose θ and the privacy-utility tradeoffs these algorithms incur. To clarify the notation

⁴Notably, Balle and Wang (2018) provide a way to achieve (ε, δ) -indistinguishability for $\varepsilon > 1$, and their technique can be adapted without loss of generality to our setting.

Algorithm 2 $(\epsilon, \delta, \theta)$ -Certified Uniformity with Exact Inverse Hessian

Require: Dataset \mathcal{D} ; forget set \mathcal{D}_f ; pretrained model $\mathbf{w}^* = \mathcal{A}(\mathcal{D})$; privacy budgets ϵ and δ ; uniformity-utility tradeoff coefficient θ ; local convex coefficient λ ; norm upper bound C .
 $\tilde{\mathbf{w}} \leftarrow \mathbf{w}^* - (\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} \nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}$. // Derived in Appx. G.
 Compute Δ as the bound in Eq. (G.20).
 $\sigma = \frac{\Delta}{\epsilon} \sqrt{2 \ln(1.25/\delta)}$.
 $\mathbf{w}^- \leftarrow \tilde{\mathbf{w}} + Y$ where $Y \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$.
return \mathbf{w}^- .

used in this section, we include a symbol table in Appx. N. Firstly, we seek to understand how we can choose θ to guarantee uniformity over the forget set. To do so, we provide a constraint to be satisfied to ensure uniformity. Then, we provide an appropriate lower bound on θ to ensure the constraint is satisfied. In doing so, one obtains a bound on the privacy of our algorithm. We next want to obtain a bound on the utility of our algorithm. To that end, we upper bound the difference between the retain loss of the locally optimal learned model $\mathcal{A}(\mathcal{D}_r)$ and the locally optimal solution to the Pareto objective $\mathcal{M}_\theta(\mathcal{D})$. We obtain a tight, non-vacuous bound with respect to θ and characterize it asymptotically. Incorporating the two bounds provides a concrete characterization of the privacy-utility tradeoffs that occur when our algorithms are used.

In particular, across all algorithms, we make the pretrained model indistinguishable from:

$$\mathcal{M}_\theta(\mathcal{D}) = \arg \min_{\|\mathbf{w}\|_2 \leq C, \mathbf{w} \in \mathcal{W}} \theta \mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (5)$$

$$= \arg \min_{\|\mathbf{w}\|_2 \leq C, \mathbf{w} \in \mathcal{W}} \theta \sum_{i=1}^{|\mathcal{D}_f|} \ell_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f^{(i)}) + (1 - \theta) \sum_{j=1}^{|\mathcal{D}_r|} \ell_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r^{(j)}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2, \quad (6)$$

where $\ell_{\mathcal{K}}$ and $\ell_{\mathcal{A}}$ are component loss functions corresponding to individual data instances in the forget and retain sets, respectively. Note that λ is present either as weight decay in the Pareto learning in Alg. 1 or as part of the local convex approximation in Alg. 2. Furthermore, note that the objective is constrained by $\|\mathbf{w}\|_2 \leq C$; we use this as a part of our local convex approximation when deriving Alg. 2; it is however unnecessary for Alg. 1. Similarly, unlearning methods assume this either implicitly or explicitly (Zhang et al., 2024). One can use projected gradient descent (Nocedal and Wright, 1999) during pretraining to satisfy this constraint.

Note that Alg. 1 has $\lambda \approx 0$. For Alg. 2, by Lemma H.9, our models $f_{\mathbf{w}^-}$ and $f_{\mathcal{M}_\theta(\mathcal{D})}$ have approximately the same outputs over \mathcal{D}_f , where \mathbf{w}^- are the weights after applying one of our TTP algorithms. Hence, for any of our algorithms, to ensure indistinguishability from uniformity over \mathcal{D}_f , it suffices to ensure that \mathcal{M}_θ satisfies the following constraint:

$$\|f_{\mathcal{M}_\theta(\mathcal{D})}(\mathcal{D}_f) - U[0, |\mathcal{Y}]\|_\infty \leq \gamma. \quad (7)$$

We then have the following bound on Eq. (7) with respect to θ , the proof of which is in Appx. H.10:

Proposition 5.1. *Let $\mathcal{M}_\theta(\mathcal{D})$ be the global solution to the Pareto objective. Choose, as surrogate losses, $\ell_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f^{(i)}) = D_{KL}(f_{\mathbf{w}}(\mathcal{D}_f^{(i)}) \| U[0, |\mathcal{Y}|])$, the KL divergence between the model outputs over the forget set and the uniform distribution, and $\ell_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r^{(j)}) = \mathbb{H}_{CE}(\mathcal{D}_r^{(j, \mathcal{Y})}, f_{\mathbf{w}}(\mathcal{D}_r^{(j, \mathcal{X})}))$, the cross entropy between model predictions and labels over the retain set. Then, $\|f_{\mathcal{M}_\theta(\mathcal{D})}(\mathcal{D}_f) - U[0, |\mathcal{Y}]\|_\infty \leq \sqrt{2(\frac{1-\theta}{\theta} |\mathcal{D}_r| \ln |\mathcal{Y}|)}$.*

Proof Sketch: By using Prop. 4.2 and the fact that $\mathcal{M}_\theta(\mathcal{D})$ is a global minimizer, we can yield a bound on $\mathcal{L}_{\mathcal{K}}$. Then, standard inequalities yield our result.

Then, we can choose θ as follows to guarantee Eq. (7), the proof of which is in Appx. H.11:

Corollary 5.2. *For a $\gamma > 0$, choosing $\theta \geq \frac{2|\mathcal{D}_r| \ln |\mathcal{Y}|}{\gamma^2 + 2|\mathcal{D}_r| \ln |\mathcal{Y}|}$ guarantees that Eq. (7) holds.*

Discussion: Note that Cor. 5.2 is well-defined in that $\theta \in (0, 1)$ for any choice of $|\mathcal{D}_r|, |\mathcal{Y}|$ and ϵ . Furthermore, Cor. 5.2 restricts $\mathcal{M}_\theta(\mathcal{D})$ to a subset of Pareto optimal solutions, but this does not

render it no longer Pareto optimal; thus, our formulation as in Appx. G still holds in its entirety. Importantly, this is a sufficient *but not necessary* condition to satisfy Eq. (7).

Similarly, by Lemma H.9, we can study the affect of θ in $\mathcal{M}_\theta(\mathcal{D})$ on the (empirical) retain error on \mathcal{D}_r , after our algorithms are applied. To provide this bound, we require two key assumptions:

Assumption 5.3. *The gradients of $\ell_{\mathcal{K}}$ and $\ell_{\mathcal{A}}$ are Lipschitz in \mathbf{w} with constants $\frac{P_{\mathcal{K}}}{|\mathcal{D}_f|}$ and $\frac{P_{\mathcal{A}}}{|\mathcal{D}_r|}$.*

Assumption 5.4. *The Hessians of $\ell_{\mathcal{K}}$ and $\ell_{\mathcal{A}}$ are Lipschitz in \mathbf{w} with constants $\frac{F_{\mathcal{K}}}{|\mathcal{D}_f|}$ and $\frac{F_{\mathcal{A}}}{|\mathcal{D}_r|}$.*

Discussion: Note that these assumptions are only used to prove Thm. 5.5 and in Appx. G; we do not require them to prove all previously mentioned theorems. These assumptions, similar to those studied by Zhang et al. (2024), are less restrictive than those typically studied in certified unlearning Sekhari et al. (2021); importantly, we do not assume (strong) convexity of the losses.

We then present a tight, non-vacuous bound on the retain error after applying any of our algorithms:

Theorem 5.5. *Suppose Assumptions 5.3 and 5.4 hold, and let $P_{\mathcal{K}}, P_{\mathcal{A}}, F_{\mathcal{K}}, F_{\mathcal{A}}$ be as defined in Assumptions 5.3 and 5.4. Let $\alpha^* := \mathcal{L}_{\mathcal{A}}(\mathcal{A}(\mathcal{D}_r), \mathcal{D}_r)$ be the locally optimal (empirical) retain loss, achieved by $\mathcal{M}_\theta(\mathcal{D})$ when $\theta = 0$. Let $\alpha(\theta) := \mathcal{L}_{\mathcal{A}}(\mathcal{M}_\theta(\mathcal{D}), \mathcal{D}_r)$ be the locally optimal retain loss obtained by $\mathcal{M}_\theta(\mathcal{D})$ when $\theta \in (0, 1)$. Suppose all weights used throughout are bounded by $\|\mathbf{w}\|_2 \leq C$. Additionally, denote by $F := \theta F_{\mathcal{K}} + (1 - \theta)F_{\mathcal{A}}$ and $P := \theta P_{\mathcal{K}} + (1 - \theta)P_{\mathcal{A}}$. Consider regularization coefficient $\lambda \geq L + 2\theta CF + \sqrt{2\theta CF(P + 2\theta CF + 8P_{\mathcal{K}})}$. Then, we have the following bound:*

$$|\alpha^* - \alpha(\theta)| \leq \mathcal{O}(\lambda C^2\theta + C^2\theta^2). \quad (8)$$

Proof Sketch: We subtract the first order conditions, by definition of α^* and $\alpha(\theta)$, to get an expression with respect to the gradients; plugging in an equivalent path integral expression and applying Lemma H.2 yields our desired result, with a full proof (including the full bound) in Appx. H.12.

Discussion: Three key hyperparameters should be kept small to ensure high retain accuracy: the ℓ_2 regularization coefficient λ , the max model weight magnitude C , and the Pareto frontier hyperparameter θ . In particular, large regularization coefficients take the model off the Pareto frontier. However, smaller or sparser weights are preferred, since the bound grows quadratically in C .

In addition, that when $\theta = 0$, the bound simplifies to 0, indicating that it is tight near 0. We demonstrate that it is tight near 1 in Appx. L. Furthermore, in the case of Alg. 1, since $\lambda \approx 0$, we do not need the condition on λ and obtain a clearer characterization. Furthermore, we can obtain a more concise bound with simpler techniques, but such a bound is vacuous and does not incorporate information about θ ; we elaborate on this in Appx. H.12.

6 EMPIRICAL ANALYSIS

Below, we provide empirical results for Alg. 1 and Alg. 2. We firstly discuss our experimental setup, baselines, and define our uniformity metric. Next, we provide our core results across Alg. 1, Alg. 2, and our baselines for several architectures and benchmarks. We also comment on the Pareto frontier of Alg. 1 and Alg. 2, providing additional insight into the structure of our problem.

Setup and Baselines. Our primary results on Alg. 1 are for ResNet50 (He et al., 2016) trained on SVHN, CIFAR10, and CIFAR100. We also provide results for logistic regression on MNIST to evaluate Alg. 2. We then include additional experiments with more complicated datasets and models, such as ViT (Dosovitskiy et al., 2021) and TinyImageNet (Le and Yang, 2015), in Tab. L.3. We compare results with the pretrained model and the model retrained without the forget set, which constitutes exact unlearning (Bourtoule et al., 2021). We also compare our methods to LabelDP (Ghazi et al., 2021) and a synthetic baseline that assigns random labels to instances neighboring the forget set. Across methods, we compare retain accuracy, test accuracy, and forget uniformity. We provide more details and the rationale for our baselines in Appx. K.

Providing a Uniformity Metric: We require a metric to compare uniformity over the forget set in an interpretable manner. Thus, we define the ‘‘confidence distance’’ as $\max\{0, f(\mathbf{x})_t - \frac{1}{|\mathcal{D}_f|}\}$ for $\mathbf{x} \in \mathcal{D}_f$, where $f(\mathbf{x})_t$ is the max confidence score. In our experiments, we use this as the primary

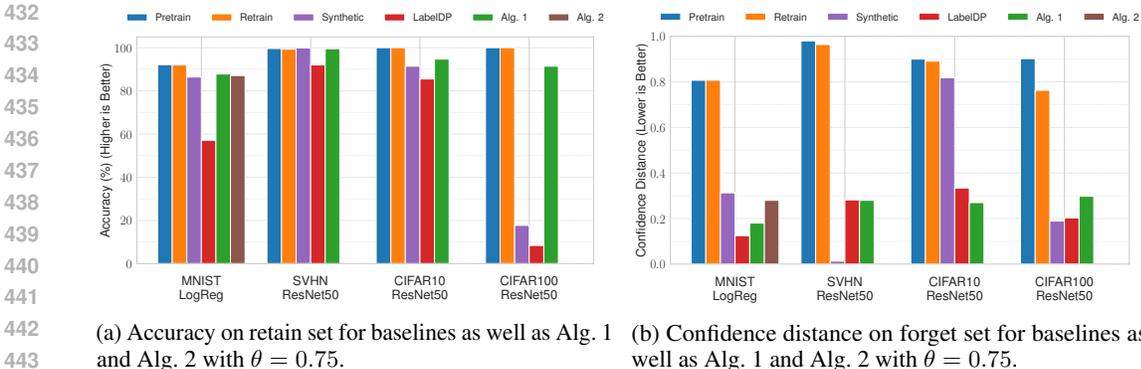


Figure 2: Across datasets, observe a significant drop in confidence distance, where lower is better, for both our algorithms. We also observe that both algorithms provide strong accuracy on the retain set. We observe similar behavior for the test set in Appx. L, while the baselines are inconsistent. Variance is negligible for all metrics.

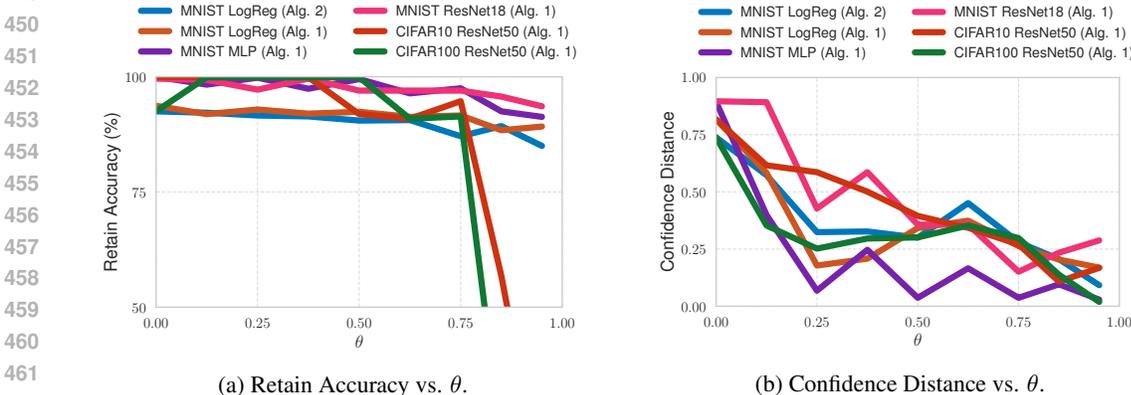


Figure 3: From Fig. 3a, we observe that for simple datasets, the retain accuracy decreases smoothly. However, for larger datasets like CIFAR10 and CIFAR100 as one passes $\theta \approx 0.75$, retain accuracy drops significantly. This motivates our choice of $\theta = 0.75$ used throughout our experiments. In Fig. 3b we observe that the confidence distance decreases roughly linearly as θ increases.

metric for uniformity, reporting the average confidence distance over the forget set. We discuss why this is reasonable in Appx. D and compare it to alternative metrics in Appx. L.

Overall Results: The results for Alg. 1 are presented in Fig. 2, in which we were able to achieve a $> 3\times$ decrease in confidence distance with only a 0.01% and 0.04% decrease in retain and test accuracy, respectively, for a ResNet50 pretrained on SVHN. We obtain similar results for MNIST, CIFAR10, and CIFAR100: retain and test set accuracies remain high, while forget confidence distance is significantly reduced. Results for the test set are deferred to Appx. L. We additionally find that the synthetic baseline can induce uniformity well for SVHN, but can either fail to induce uncertainty entirely (CIFAR10) or induce uncertainty at great cost to retain and test accuracy (CIFAR100). We observe similar behavior for TinyImageNet in Appx. L.3. This holds similarly for LabelDP, which furthermore undesirably reduces the confidence distance on retain and test sets, while our method does not, as demonstrated in Tab. L.4. Furthermore, our observations coincide with Zhao et al. (2024), observing that unlearned models still produce confident predictions on deleted instances.

Furthermore, as illustrated by Fig. 2, we find that Alg. 2 also induces uniformity well, while marginally reducing retain and test accuracy. Thus, this algorithm produces a certificate through which test-time privacy can be verified while still obtaining a good privacy-utility tradeoff. For both algorithms, tables are included in Appx. L for completeness.

Pareto Frontiers: To better understand the structure of our problem, we explore the Pareto frontier in Fig. 3. We observe that for MNIST, CIFAR10, and CIFAR100, various θ can provide good retain

accuracy, albeit at the cost of uniformity. In general, we find that $\theta \approx 0.75$ offers a solid privacy-utility tradeoff. Thus, the ε in Prop. 5.1 can be chosen fairly large while ensuring low confidence distance.

Additional Experiments: We conduct various additional experiments in in Appx. L and briefly comment about them here. Firstly, we obtain excellent performance for TinyImageNet and ViT in Appx. L.3. Secondly, as desired, we obtain high confidence distances on the retain and test sets in Appx. L.4. Thirdly, we study the optimization dynamics of Alg. 1 in Appx. L.6, providing mathematical and empirical evidence for the necessity of early stopping in large models when using Alg. 1. Fourthly, we evaluate our method on several strong TTP attacks, demonstrating that we can still offer effective defense, especially when compared to pretraining or retraining, in Appx. L.7. Fifthly, in Appx. L.8, we find that we preserve strong accuracy and high confidence, as desired, on test instances which are nearest neighbors to the forget set instances. Thus, an adversary querying nearby instances outside of the forget set does not suffice to circumvent our algorithms. Sixthly, we find that we can induce uncertainty on forget instances which were not part of the original training dataset, while still preserving retain and test accuracies, in Appx. L.9. Seventhly, we provide ablations on the size of our forget set in Appx. L.10. Finally, we compare our confidence metric to an ℓ_2 uniformity metric, finding that they highly correlate, in Appx. L.11.

7 DISCUSSION

We present *test-time privacy*, a threat model in which an adversary seeks to directly use a confident prediction for harm. This contrasts with existing work like PATE and LabelDP, which focus on protecting against model inversion and leakage of ground truth labels. To protect against a test-time privacy adversary, we present multiple algorithms to induce uniformity on a known corrupted subset while preserving utility on the rest of the data instances. This can be used to prevent adversaries from taking advantage of model outputs. Furthermore, we prove a privacy-utility tradeoff for our algorithms, providing a tight bound which is empirically verified. We hope our test-time privacy can further inspire the community to explore different threat models for sensitive data. Limitations and future directions are provided in Appx. F.

REFERENCES

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. *Proceedings of the Conference on Computer and Communications Security (CCS)*, pages 308–318, 2016.
- Naman Agarwal, Brian Bullins, and Elad Hazan. Second-order stochastic optimization in linear time. *stat*, 1050:15, 2016.
- Mohammad Al-Rubaie and J Morris Chang. Privacy-preserving machine learning: threats and solutions. *IEEE Security & Privacy*, 17(2):49–58, 2019.
- Youssef Allouah, Joshua Kazdan, Rachid Guerraoui, and Sanmi Koyejo. The utility and complexity of in-and out-of-distribution machine unlearning. *International Conference on Machine Learning (ICML)*, 2025.
- Kareem Amin, Alex Kulesza, and Sergei Vassilvitskii. Practical considerations for differential privacy, 2024. URL <https://arxiv.org/abs/2408.07614>.
- Anastasios N Angelopoulos, Michael I Jordan, and Ryan J Tibshirani. Gradient equilibrium in online learning: theory and applications. *arXiv preprint arXiv:2501.08330*, 2025.
- Yoshinori Aono, Takuya Hayashi, Lihua Wang, Shiho Moriai, et al. Privacy-preserving deep learning via additively homomorphic encryption. *IEEE Transactions on Information Forensics and Security*, 13(5):1333–1345, 2017.
- Mirza Ahad Baig and Krzysztof Pietrzak. On the (in)security of proofs-of-space based longest-chain blockchains. Cryptology ePrint Archive, Paper 2025/942, 2025. URL <https://eprint.iacr.org/2025/942>.

- 540 Borja Balle and Yu-Xiang Wang. Improving the gaussian mechanism for differential privacy:
541 analytical calibration and optimal denoising. *International Conference on Machine Learning*
542 (*ICML*), 2018.
- 543
544 Lucas Bourtole, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers,
545 Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. *IEEE Symposium on Security*
546 *and Privacy (SP)*, pages 141–159, 2021.
- 547 Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. (leveled) fully homomorphic encryption
548 without bootstrapping. *Transactions on Computation Theory (TOCT)*, 6(3):1–36, 2014.
- 549
550 Mark Bun, Damien Desfontaines, Cynthia Dwork, Naor Moni, Kobbi Nissim, Aaron Roth, Adam
551 Smith, Thomas Steinke, Jonathan Ullman, and Salil Vadhan. Statistical inference is not a privacy
552 violation, 2021.
- 553 Robert Istvan Busa-Fekete, Umar Syed, Sergei Vassilvitskii, et al. On the pitfalls of label differential
554 privacy. In *NeurIPS Workshops*, 2021.
- 555
556 Sungmin Cha, Sungjun Cho, Dasol Hwang, Honglak Lee, Taesup Moon, and Moontae Lee. Learning
557 to unlearn: instance-wise unlearning for pre-trained classifiers. *AAAI Conference on Artificial*
558 *Intelligence*, 38(10):11186–11194, 2024.
- 559 Kamalika Chaudhuri and Claire Monteleoni. Privacy-preserving logistic regression. *Advances in*
560 *Neural Information Processing Systems (NeurIPS)*, 21, 2008.
- 561
562 Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical risk
563 minimization. *Journal of Machine Learning Research*, 12(3), 2011.
- 564
565 Lynn Chua, Badih Ghazi, Pritish Kamath, Ravi Kumar, Pasin Manurangsi, Amer Sinha, and Chiyuan
566 Zhang. Scalable dp-sgd: shuffling vs. poisson subsampling. *Advances in Neural Information*
567 *Processing Systems (NeurIPS)*, 37:70026–70047, 2024.
- 568
569 Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David
570 Ha. Deep learning for classical japanese literature. *Advances in Neural Information Processing*
Systems (NeurIPS), 2018.
- 571
572 Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal*
Processing, 29(6):141–142, 2012.
- 573
574 Irit Dinur and Kobbi Nissim. Revealing information while preserving privacy. *Proceedings of the*
575 *Symposium on Principles of Databases (PODS)*, pages 202–210, 2003.
- 576
577 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
578 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
579 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.
International Conference on Learning Representations (ICLR), 2021.
- 580
581 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in
582 private data analysis. *Proceedings of the Theory of Cryptography Conference (TCC)*, 2006.
- 583
584 Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations*
585 *and Trends® in Theoretical Computer Science*, 9(3-4):211–407, 2014.
- 586
587 European Parliament. Regulation (EU) 2016/679 of the European Parliament and of the Council,
2016. URL <https://data.europa.eu/eli/reg/2016/679/oj>.
- 588
589 Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence
590 information and basic countermeasures. In *Proceedings of the Conference on Computer and*
591 *Communications Security (CCS)*, pages 1322–1333, 2015.
- 592
593 Badih Ghazi, Noah Golowich, Ravi Kumar, Pasin Manurangsi, and Chiyuan Zhang. Deep learning
with label differential privacy. *Advances in Neural Information Processing Systems (NeurIPS)*, 34:
27131–27145, 2021.

- 594 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
595 examples. *International Conference on Learning Representations (ICLR)*, 2015.
- 596
- 597 Google. Vision transformer pretrained on imagenet-21k, 2023. URL <https://huggingface.co/google/vit-base-patch16-224>.
- 598
- 599 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
600 networks. *International Conference on Machine Learning (ICML)*, 2017.
- 601
- 602 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
603 recognition. *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778,
604 2016.
- 605 C-L Hwang and Abu Syed Md Masud. *Multiple objective decision making—methods and applications:
606 a state-of-the-art survey*, volume 164. Springer Science & Business Media, 2012.
- 607
- 608 Gautam Kamath. Lecture 12: what is privacy? *Lectures on private ml and stats*, 2020.
- 609
- 610 Anastasia Koloskova, Youssef Allouah, Animesh Jha, Rachid Guerraoui, and Sanmi Koyejo. Certified
611 unlearning for neural networks. *International Conference on Machine Learning (ICML)*, 2025.
- 612 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
613 *Technical Report from the University of Toronto*, 2009.
- 614
- 615 Solomon Kullback and Richard A Leibler. On information and sufficiency. *The Anals of Mathematical
616 Statistics*, 22(1):79–86, 1951.
- 617 Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. Towards unbounded
618 machine unlearning. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- 619
- 620 Ya Le and Xuan S. Yang. Tiny imagenet visual recognition challenge, 2015.
- 621
- 622 Joon-Woo Lee, HyungChul Kang, Yongwoo Lee, Woosuk Choi, Jieun Eom, Maxim Deryabin,
623 Eunsang Lee, Junghyun Lee, Donghoon Yoo, Young-Sik Kim, et al. Privacy-preserving machine
624 learning with fully homomorphic encryption for deep neural network. *IEEE Access*, 10:30039–
30054, 2022.
- 625
- 626 Ninghui Li, Wahbeh Qardaji, and Dong Su. On sampling, anonymization, and differential privacy or,
627 k-anonymization meets differential privacy. *Proceedings of the Conference on on Computer and
628 Communications Security (CCS)*, pages 32–33, 2012.
- 629
- 630 Hui Liu, Yibo Dou, Kai Wang, Yunmin Zou, Gan Sen, Xiangtao Liu, and Huling Li. A skin disease
631 classification model based on multi scale combined efficient channel attention module. *Scientific
632 Reports*, 15(1):6116, 2025.
- 633
- 634 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
635 Towards deep learning models resistant to adversarial attacks. *International Conference on
636 Learning Representations (ICLR)*, 2018.
- 637
- 638 Günter Mayer. On the convergence of the neumann series in interval analysis. *Linear algebra and its
639 applications*, 65:63–70, 1985.
- 640
- 641 Frank Mcsherry. Statistical inference considered harmful, 2016. URL [https://github.com/
642 frankmcsherry/blog/blob/master/posts/2016-06-14.md](https://github.com/frankmcsherry/blog/blob/master/posts/2016-06-14.md).
- 643
- 644 Ronak Mehta, Sourav Pal, Vikas Singh, and Sathya N Ravi. Deep unlearning via randomized
645 conditionally independent Hessians. *Conference on Computer Vision and Pattern Recognition
646 (CVPR)*, 2022.
- 647
- 648 Merriam-Webster. Privacy. *Merriam-Webster Dictionary*, 2022. URL [https://www.
649 merriam-webster.com/dictionary/privacy](https://www.merriam-webster.com/dictionary/privacy).
- 650
- 651 Microsoft. Resnet50 pretrained on imagenet-21k, 2024. URL [https://huggingface.co/
652 microsoft/resnet-50](https://huggingface.co/microsoft/resnet-50).

- 648 Kaisa Miettinen. *Nonlinear multiobjective optimization*, volume 12. Springer Science & Business
649 Media, 1999.
- 650
- 651 Siqiao Mu and Diego Klabjan. Rewind-to-delete: certified machine unlearning for nonconvex
652 functions, 2025. URL <https://arxiv.org/abs/2409.09778>.
- 653
- 654 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al.
655 Reading digits in natural images with unsupervised feature learning. *NeurIPS Workshops*, page 4,
656 2011.
- 657 Thanh Tam Nguyen, Thanh Trung Huynh, Zhao Ren, Phi Le Nguyen, Alan Wee-Chung Liew,
658 Hongzhi Yin, and Quoc Viet Hung Nguyen. A survey of machine unlearning. *arXiv preprint*
659 *arXiv:2209.02299*, 2022.
- 660
- 661 Jorge Nocedal and Stephen J Wright. *Numerical optimization*. Springer, 1999.
- 662 International Association of Privacy Professionals. Swedish court rejects
663 google’s appeal in rtbf case, 2020. URL [https://iapp.org/news/a/
664 swedish-court-rejects-googles-appeal-in-rtbf-case](https://iapp.org/news/a/swedish-court-rejects-googles-appeal-in-rtbf-case).
- 665
- 666 Nicolas Papernot, Shuang Song, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and Úlfar Erlings-
667 son. Scalable private learning with pate. *International Conference on Learning Representations*
668 (*ICLR*), 2018.
- 669 Panos M Pardalos, Antanas Žilinskas, Julius Žilinskas, et al. *Non-convex multi-objective optimization*.
670 Springer, 2017.
- 671
- 672 Tim Pearce, Alexandra Brintrup, and Jun Zhu. Understanding softmax confidence and uncertainty.
673 *arXiv preprint arXiv:2106.04972*, 2021.
- 674
- 675 Barak A Pearlmutter. Fast exact multiplication by the hessian. *Neural computation*, 6(1):147–160,
676 1994.
- 677 Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz Kaiser, and Geoffrey Hinton. Regularizing
678 neural networks by penalizing confident output distributions. *International Conference on Learning*
679 *Representations (ICLR)*, 2017.
- 680
- 681 Mark S Pinsker. Information and information stability of random variables and processes. *Holden-*
682 *Day*, 1964.
- 683
- 684 Xinbao Qiao, Meng Zhang, Ming Tang, and Ermin Wei. Hessian-free online certified unlearning.
685 *International Conference on Learning Representations (ICLR)*, 2025.
- 686
- 687 Stefan Schoepf, Michael Curtis Mozer, Nicole Elyse Mitchell, Alexandra Brintrup, Georgios Kaissis,
688 Peter Kairouz, and Eleni Triantafillou. Redirection for erasing memory (rem): Towards a universal
unlearning method for corrupted data. *arXiv preprint arXiv:2505.17730*, 2025.
- 689
- 690 Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. Remember what
691 you want to forget: algorithms for machine unlearning. *Advances in Neural Information Processing*
692 *Systems (NeurIPS)*, 34:18075–18086, 2021.
- 693
- 694 Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks
695 against machine learning models. *IEEE Symposium on Security and Privacy (SP)*, pages 3–18,
2017.
- 696
- 697 Congzheng Song, Thomas Ristenpart, and Vitaly Shmatikov. Machine learning models that remember
698 too much. *Proceedings of the Conference on on Computer and Communications Security (CCS)*,
699 pages 587–601, 2017.
- 700
- 701 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
Poole. Score-based generative modeling through stochastic differential equations. *International*
Conference on Learning Representations (ICLR), 2021.

- 702 Jingwei Sun, Ang Li, Binghui Wang, Huanrui Yang, Hai Li, and Yiran Chen. Soteria: Provable
703 defense against privacy leakage in federated learning from representation perspective. *Conference*
704 *on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- 705 Xiaoxiao Sun, Jufeng Yang, Ming Sun, and Kai Wang. A benchmark for automatic visual classification
706 of clinical skin disease images. *European Conference on Computer Vision (ECCV)*, 2016.
- 707 TorchVision. Torchvision: Pytorch’s computer vision library, 2016. URL <https://github.com/pytorch/vision>.
- 708 Cuong Tran and Ferdinando Fioretto. Personalized privacy auditing and optimization at test time.
709 *arXiv preprint arXiv:2302.00077*, 2023.
- 710 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
711 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing*
712 *Systems (NeurIPS)*, 30, 2017.
- 713 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*,
714 volume 47. Cambridge university press, 2018.
- 715 Zhibo Wang, Mengkai Song, Zhifei Zhang, Yang Song, Qian Wang, and Hairong Qi. Beyond inferring
716 class representatives: User-level privacy leakage from federated learning. *IEEE Conference on*
717 *Computer Communications*, 2019.
- 718 Annie White. Dmvs can (and do) collect and sell your personal data. *Car and Driver*, 2020.
- 719 Ruihan Wu, Jin Peng Zhou, Kilian Q Weinberger, and Chuan Guo. Does label differential privacy
720 prevent label inference attacks? *International Conference on Artificial Intelligence and Statistics*
721 *(AISTATS)*, 2023.
- 722 Taihong Xiao, Yi-Hsuan Tsai, Kihyuk Sohn, Manmohan Chandraker, and Ming-Hsuan Yang. Ad-
723 versarial learning of privacy-preserving and task-oriented representations. *AAAI Conference on*
724 *Artificial Intelligence*, 2020.
- 725 Jufeng Yang, Xiaoxiao Sun, Jie Liang, and Paul L Rosin. Clinical skin lesion diagnosis using
726 representations inspired by dermatologist criteria. *Conference on Computer Vision and Pattern*
727 *Recognition (CVPR)*, 2018.
- 728 Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning:
729 Analyzing the connection to overfitting. *2018 IEEE Symposium on Computer Security Foundations*
730 *(CSF)*, pages 268–282, 2018.
- 731 Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath, Janardhan
732 Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, et al. Differentially private fine-tuning of
733 language models. *International Conference on Learning Representations (ICLR)*, 2022.
- 734 Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn.
735 Gradient surgery for multi-task learning. *Advances in Neural Information Processing Systems*
736 *(NeurIPS)*, 33:5824–5836, 2020.
- 737 Binchi Zhang, Yushun Dong, Tianhao Wang, and Jundong Li. Towards certified unlearning for deep
738 neural networks. *International Conference on Machine Learning (ICML)*, 2024.
- 739 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep
740 learning requires rethinking generalization. *International Conference on Learning Representations*
741 *(ICLR)*, 2017.
- 742 Rui Zhang, Song Guo, Junxiao Wang, Xin Xie, and Dacheng Tao. A survey on gradient inversion:
743 Attacks, defenses and future directions. *International Joint Conferences on Artificial Intelligence*
744 *(IJCAI)*, 2022.
- 745 Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. idlg: Improved deep leakage from gradients.
746 *arXiv preprint arXiv:2001.02610*, 2020.

756 Kairan Zhao, Meghdad Kurmanji, George-Octavian Bărbulescu, Eleni Triantafillou, and Peter Tri-
757 antafillou. What makes unlearning hard and what to do about it. *Advances in Neural Information*
758 *Processing Systems (NeurIPS)*, 2024.

759
760 Mingyi Zhou, Xiang Gao, Jing Wu, John Grundy, Xiao Chen, Chunyang Chen, and Li Li. Modelob-
761 fuscator: obfuscating model information to protect deployed ml-based systems. *Proceedings of the*
762 *Symposium on Software Testing and Analysis*, pages 1005–1017, 2023.

763 Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. *Advances in Neural*
764 *Information Processing Systems (NeurIPS)*, 2019.

765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

Appendix

Table of Contents

A Full Test-Time Privacy Threat Model	18
B Defining Test-Time Privacy Attacks	19
C Additional Related Work	20
D Uniformity Metric	21
E Test-Time Privacy Examples	21
F Limitations and Future Directions	22
G Designing Certified Algorithms	23
H Proofs	26
H.1 Helpful Lemmas	26
H.2 Proof of Proposition 4.2	30
H.3 Proof of Proposition 4.3	30
H.4 Proof of Theorem G.1	31
H.5 Proof of Proposition G.2	31
H.6 Proof of Proposition G.3	32
H.7 Proof of Theorem G.4	32
H.8 Proof of Theorem G.5	32
H.9 Proof of Proposition I.1	34
H.10 Proof of Proposition 5.1	34
H.11 Proof of Corollary 5.2	35
H.12 Proof of Theorem 5.5	36
I Online Algorithm	40
J Eliminating Hyperparameters in Certified Algorithms	40
K Experimental Details	41
K.1 Dataset Details	41
K.2 Model Details	41
K.3 Baseline Details	42
K.4 Hyperparameter Details	42
L Additional Experiments	45
L.1 Test Set Accuracies for Main Paper Experiments	45
L.2 Tables for Main Paper Experiments	45
L.3 Additional Experiments on TinyImageNet & ViT	46
L.4 LabelDP and Alg. 1 Confidence Distances for Retain, Test, and Forget Sets	47
L.5 Pareto Frontier Main Paper Table	47
L.6 Optimization Dynamics	48
L.7 Evaluating Test-Time Privacy Attacks	51
L.8 Robustness of Alg. 1 Classifier on Neighboring Test Instances	52
L.9 Ensuring Test-Time Privacy for Test Instances	52
L.10 Ablation Study on Forget Set Size	53
L.11 Evaluating Confidence Distance as a TTP Metric	54
L.12 An Additional Baseline with Randomly Sampled Labels: GaussianUniform	56

864	L.13 Tightness of Bound in Theorem 5.5	57
865	L.14 Confidence Intervals for Main Paper Experiments	57
866	L.15 Proportions of Time Elapsed in Alg. 1	57
867	L.16 Warmup Values for MNIST LogReg	57
868	L.17 Visualization of Softmax Outputs	58
869	L.18 Results for KMNIST LogReg and MLP	60
870	L.19 Ablation Study on Synthetic Baseline Sample Size	60
871	L.20 Test Accuracy Plot for Pareto Frontier Experiments	61
872	L.21 Results for Alg. 3	61
873	L.22 Hyperparameter Sensitivity for Alg. 3	61
874		
875	M Broader Impacts	63
876		
877	N Symbol Table	64
878		
879		
880		
881		
882		
883		
884		
885		
886		
887		
888		
889		
890		
891		
892		
893		
894		
895		
896		
897		
898		
899		
900		
901		
902		
903		
904		
905		
906		
907		
908		
909		
910		
911		
912		
913		
914		
915		
916		
917		

A FULL TEST-TIME PRIVACY THREAT MODEL

Following recent works on privacy and cybersecurity (Baig and Pietrzak, 2025), we begin by making our threat model concrete as an informal security game. Broadly, we consider a *test-time privacy (TTP) game* where a *TTP adversary* aims use an open-weight ML model \hat{f} to produce a confident, harmful prediction m for a specific set of corrupted inputs \mathcal{D}_f drawn from a distribution \mathcal{P} .

Actors and Assets: The game begins with three key actors:

1. The data corrupter c , an entity that either maliciously or erroneously creates a “forget set” \mathcal{D}_f of corrupted instances, e.g. a server which makes an error in compressing a medical image uploaded to an online forum.
2. The *model provider* τ , a benign challenger that uses a learning algorithm \mathcal{A} and releases a model f . For example, f can classify skin disease from skin images. They then seek to ensure TTP by running algorithm \mathcal{G} to obtain model \hat{f} .
3. The TTP adversary ν e.g. a potential medical insurance provider who has access to the architecture and parameters of \hat{f} and aims to obtain harmful prediction m on \mathcal{D}_f to e.g. use as a warrant to reject insurance applicants.

We operate under two core assumptions:

Assumption A.1 (Open-Weight Access). *The adversary ν has complete access to the model’s architecture and weights.*

Assumption A.2. τ -Limited Knowledge *The model provider τ is notified about the existence of a corrupted forget set \mathcal{D}_f , but does not know the specific harmful label m . Furthermore, they do not know the specific adversary ν .*

Remarks: Specifically, Asm. A.1 renders naive defenses, like obtaining \hat{f} by masking softmax outputs of f , useless, as an adversary can simply remove such a mask and recover prediction m . Furthermore, to make Asm. A.2 concrete, following the example provided in Sec. 1, the model provider τ does not know whether e.g. ν is a medical insurance company aiming to obtain a prediction of “Melanoma” to reject coverage or a defense attorney in a criminal case against a doctor aiming to obtain a prediction of “Benign” to clear a doctor of accusations of medical malpractice.

Game: The game is then played in the first round.

Round 1: The first round contains preliminary steps as follows:

- The corrupter c corrupts the data and yields \mathcal{D}_f , which adversary ν gains access to e.g. through the public Internet.
- The model provider τ trains a model f over instances from \mathcal{P} .
- The model provider τ is made aware that \mathcal{D}_f contains corrupted instances, and seeks to protect them from a TTP adversary ν .

Round 2: The second round contains the following steps:

1. The model provider τ , who is aware of TTP, aims to provide a model \hat{f} to replace f such that $\hat{f}(x) \neq m$, where m is the harmful prediction. However, they are *unaware of which prediction m is*. τ thus runs an algorithm \mathcal{G} with respect to f , \mathcal{D}_f , and a training dataset \mathcal{D} , which yields a new model \hat{f} .
2. The TTP adversary takes model \hat{f} and attempts to obtain a confident prediction m which serves as a warrant to endanger individuals e.g. to reject individuals from a health insurance provider because their image was classified as a high risk disease like melanoma.

Win Conditions: *The TTP adversary ν wins if it can leverage the provided information about \hat{f} , including the model weights of \hat{f} , to confidently recover the sensitive prediction m for any $x \in \mathcal{D}_f$, as they can then e.g. use this prediction as a warrant to reject people’s insurance applications. The model provider τ wins if \hat{f} leaves ν uncertain as to whether the prediction is m or not.*

Given this win condition, an algorithm \mathcal{G} satisfies *test-time privacy* if the adversary can only guess at the model output for all instances in \mathcal{D}_f . Thus, it is optimal to induce maximal uncertainty over \mathcal{D}_f , since the model provider τ does not know beforehand which sensitive class results in a TTP violation. We make this concrete in Appx. C, where we discuss existing methods to misclassify and relabel training instances.

In particular, in the discriminative setting—which our work focuses on—it is optimal for model f to output uniform softmax outputs over \mathcal{D}_f , while maintaining strong accuracy on all other instances. Furthermore, to defend against such an adversary with open-weight model access, one must perturb the model weights in a non-invertible manner, motivating our approaches detailed in Sec. 4.

Importantly, while in our formulation in Sec. 4 we define the forget set of corrupted instances in terms of the training dataset, we do so **without loss of generality**. As detailed previously, we assume that the forget set contains *all* corrupted instances, including instances outside of the training dataset that are known to be corrupted. Denoting the set of training forget set instances $\mathcal{D}_f^{\text{train}}$ and $\mathcal{D}_f^{\text{test}}$, we can thus let $\mathcal{D}_f = \mathcal{D}_f^{\text{train}} \cup \mathcal{D}_f^{\text{test}}$ and again consider $\mathcal{D} = \mathcal{D}_f \cup \mathcal{D}_r$; in this scenario, all formal definitions and statements throughout Sec. 4, Sec. 5, and elsewhere follow in the exact same manner. A concrete example of when instances outside of a training dataset can become relevant is credit score classification; one’s credit score report can become corrupted, even if they are not in the training dataset, and one should be able to ask a credit bureau to remedy this to ensure that e.g. a loan officer does not incorrectly estimate their credit score.

Next, our work is motivated by *prevention* rather than *correction*. In particular, we assume that the adversary leverages the model \hat{f} as a part of an automated system by which they make decisions about people, including those in the forget set \mathcal{D}_f . The goal of TTP is to prevent this automation from causing harm *before* the harm occurs. Thus this eliminates the need for a costly, time consuming, and technically challenging appeals process.

Finally, in Appx. B, we present some simple TTP attacks on open-weight image classifiers to further motivate our threat model.

B DEFINING TEST-TIME PRIVACY ATTACKS

In what follows, in light of our threat model provided in Appx. A, we design some simple test-time privacy attacks to motivate our problem. We also include experiments on these attacks, and how Alg. 1 performs against them, in Appx. L.7.

Our first simple algorithm is to add a small amount of uniformly sampled Gaussian noise, presented in Alg. 5. We find that this is not very effective in increasing confidence distance, as demonstrated in Tab. L.8 and Tab. L.9. When it brings the confidence distances from low to moderate, the model is usually confidently wrong, as demonstrated in Tab. L.10.

One way to more optimally attack the TTP of a pretrained model is by finding instances in a ϕ -ball around the forget set instances that maximize the prediction confidence. To design such an attack, suppose we have a pretrained classifier $f_{w^*} : \mathcal{X} \rightarrow \Delta_{|\mathcal{Y}|}$. Here, $f_w(\mathbf{x}) = \text{softmax}(\mathbf{z}(\mathbf{x}))$, for $\mathbf{x} \in \mathcal{X}$, where \mathbf{z} is a vector of logits. For a forget set instance, we begin by adding a small amount of uniform noise to break symmetry and obtain a nonzero gradient. We then want to obtain the worst-case perturbation over the logits by solving the optimization problem:

$$\max_{\phi} \max_j \mathbf{z}_{w^*}^{(j)}(\mathbf{x} + \phi), \quad (\text{B.9})$$

$$\text{s.t. } \|\phi\|_{\infty} \leq \gamma. \quad (\text{B.10})$$

Since the max function is not differentiable everywhere, we use LogSumExp to approximate it. Denote $\rho(f_w(\mathbf{x} + \phi)) = \log \sum_{j=1}^{|\mathcal{Y}|} \exp(\mathbf{z}_{w^*}^j(\mathbf{x} + \phi))$. This yields the optimization problem:

$$\max_{\phi} \rho(f_{w^*}(\mathbf{x} + \delta)), \quad (\text{B.11})$$

$$\text{s.t. } \|\phi\|_{\infty} \leq \gamma. \quad (\text{B.12})$$

Following the Fast Gradient Sign Method (FGSM), a simple attack used to generate adversarial examples (Goodfellow et al., 2015), we design an attack as Alg. 6. Intuitively, we take a single linear step towards maximizing the function. We also design a stronger attack based on Projected Gradient Descent (PGD) (Madry et al., 2018) as Alg. 7, taking 40 steps while incrementally maximizing the confidence function while projecting back to the ball around the original instance. Empirical results are in Appx. L.7.

C ADDITIONAL RELATED WORK

Differential Privacy: Differential privacy has widely been studied in the ML community in order to ensure privacy-preservation (Chaudhuri and Monteleoni, 2008); (Chaudhuri et al., 2011); (Abadi et al., 2016); (Chua et al., 2024). There also exist methods to finetune pretrained models to satisfy differential privacy (Yu et al., 2022). Furthermore, there are also ways to aggregate label noise to preserve privacy (Papernot et al., 2018).

However, differential privacy is designed to address an entirely different threat model than ours. In particular, in the threat model of differential privacy, an adversary seeks to use model outputs to recover private information about data instance x_p corresponding to person p with e.g. a model inversion attack. A differentially private classifier generally results in confident, accurate predictions. This does not address our threat model, where an adversary may use confident model outputs to violate the privacy of person p in a different manner, taking advantage of them directly to use as a warrant to cause harm to person p .

Label Differential Privacy: Similarly, our formulation differs from label differential privacy (LabelDP) (Ghazi et al., 2021), which seeks to protect an adversary from learning the true labels of the instances in the training data. Given an instance, even after computing $f(x_p)$, under LabelDP an adversary cannot be confident that $f(x_p) = y$. However, LabelDP is applied to the entire dataset; our threat model involves only a particular subset of the training data. Furthermore, we do not need to protect the user’s ground truth label, necessarily. In our law enforcement example in Sec. 1, the agency does not care about the ground truth label. Instead, they want any confirmation such that they have a warrant to act adversarially towards person p ; for this, a confident prediction by model f suffices. Finally, LabelDP results in poor train and test accuracy for larger datasets e.g. CIFAR100, as demonstrated in Fig. 2.

Furthermore, from the perspective of protecting the privacy of the labels themselves, rather than protecting against any confident prediction, Busa-Fekete et al. (2021) demonstrate that testing a model, trained with LabelDP, on the training dataset allows an adversary to recover the labels of the label-private data with high probability. Since our algorithms induce uniformity, an adversary cannot infer the correct forget set labels by testing the model on the training dataset; thus, we provide better privacy against this threat model than LabelDP as well. Wu et al. (2023) argue that, under this threat model where one seeks to protect the labels, any model that generalizes must leak the accurate labels when tested on the training data. However, as we demonstrate by inducing uniformity while maintaining high test accuracy, this only holds when the model is to be tested on the *entire* training data, not a *subset* of the training data (or other test instances which are known to be corrupted), as in our setting.

Label Model Inversion Attacks: Related to LabelDP are model inversion attacks to recover the ground truth labels, like gradient inversion (Zhang et al., 2022); (Zhu et al., 2019); (Zhao et al., 2020). Yet, these methods do not report the confidence values for the recovered labels. Thus, they do not constitute test-time privacy attacks within our threat model. Furthermore, by the same token as above, an adversary seeks to recover a confident prediction to use as a warrant, not necessarily the ground truth labels. Still, these methods could potentially be extended to test-time privacy attacks by reporting a confidence score for the recovered labels. We leave this to future work.

Other Paradigms in Privacy: Other paradigms in the privacy literature correspond to a notion of “test-time privacy” which differ from our threat model. For example, several works study defense against model inversion attacks as test-time privacy (Wang et al., 2019); (Xiao et al., 2020); (Sun et al., 2021); (Tran and Fioretto, 2023). However, this is a separate threat model from ours; the adversary already has access to the instance x_p within our threat model.

Misclassification & Relabeling in Machine Unlearning: Recently, methods have emerged to finetune a model to misclassify rather than mimicking retraining from scratch (Cha et al., 2024). There are other similar relabeling methods in the debiasing literature which could be used for this purpose (Angelopoulos et al., 2025). However, these methods often achieve poor performance on the remaining training data and fail to provide protection against our threat model in all cases. In particular, a purposefully incorrect classification can also be used to endanger an individual. **To make this concrete, following the medical insurance example in Sec. 1, suppose the model provider misclassified x_p as “Benign” instead of “Melanoma” to protect against the insurance provider. But, as it turns out, the model user is a medical facility instead of the insurance provider. Thus, doctors may miss out on a critical medical diagnosis for a person p , again resulting in harm. We adjusted our existing justification for this in Appendix C to make this more clear.** Furthermore, in the binary classification case, if an adversary knows that x_p is in the forget set, they can recover the true $f(x_p)$ by taking complements, if an unlearning method which seeks to induce misclassification is used. They can also use the information that learned representations are markedly different than other similar examples to understand the method used. Additionally, in the multiclass setting, an adversary can still take complements of this class, yielding a probability of recovering the true class which is significantly better than choosing uniformly at random. Instead, it is fairer and more robust to have an output that is maximally uncertain.

Model Calibration and Confidence: In our setting, we use the model softmax outputs to represent the adversary’s confidence in the final prediction. However, some argue that this type of interpretation is incorrect, i.e. ML models are poorly calibrated (Guo et al., 2017). Still, this interpretation is common (Pearce et al., 2021), and thus a model user would likely rely on the softmax outputs as the confidence scores. We leave to inducing uncertainty over the calibrated outputs to future work.

D UNIFORMITY METRIC

The confidence distance quantifies the adversary’s confidence in their final prediction, i.e. the difference between the argmax softmax score and the uniform softmax score. Importantly, our method aims to have the adversary lack confidence in their final prediction. Thus, our metric captures what we aim to measure and is interpretable, since it is minimized at 0.

Furthermore, confidence distance allows us to quantify how uncertain the model is without relying on accuracy, since a drop in forget set accuracy is not the goal of our formulation. Next, if the maximum confidence score is very close to the uniform distribution, the probability mass of the output distribution must be distributed over the other softmax outputs, clearly yielding that the higher our uniformity metric, the more confident our model is, and the lower our uniformity metric, the less confident our model is. Additionally, it takes the dataset into account; for example, in CIFAR10, one would expect a uniformity score of ≈ 0.2 to be reasonable, as then the adversary can only be $\approx 30\%$ confident that they have a useful prediction. However, for CIFAR100, a uniformity score of ≈ 0.2 is much better, as it implies that an adversary can only be $\approx 21\%$ confident that they have a useful prediction.

One objection to the use of this metric may be that it does not indicate uniformity if it is low. For example, on CIFAR10, one could have a confidence score of 0.2, which yields that the max softmax output is 0.35. There could be three other nonzero softmax outputs of 0.3, 0.3, 0.1, 0.05; this clearly is not uniform. However, this ensures test-time privacy; a test-time privacy adversary now has little confidence in their prediction, even if they choose the first one, rendering their warrant for misuse of sensitive data useless.

We empirically compare our confidence distance metric to other similar metrics in Appx. L, finding that when our confidence metric is minimized, other metrics are minimized.

E TEST-TIME PRIVACY EXAMPLES

Here, we provide a set of examples of the TTP threat model:

Health Insurance: Suppose an open-weight medical imaging model f is released, designed to perform multiclass classification of skin photos into categories like “Dysplastic Nevus”, “Benign Keratosis”, which are usually harmless, or serious classes like “Melanoma” (Sun et al., 2016). A

1134 person p posts a photo of a harmless birthmark on his arm to a public health forum to ask a question.
 1135 During the upload, an e.g. server error or compression issue causes the image file to become corrupted,
 1136 severely distorting the birthmark. This results in a photo x_p . Next, a health insurance startup decides
 1137 to build risk profiles by scraping these public forums. They download the open-weight model f to
 1138 automatically screen images for potential health liabilities. When they feed x_p into f , it confidently
 1139 classifies x_p as “Melanoma”. This erroneous classification is then automatically added to person p ’s
 1140 risk profile, resulting person p being unfairly denied coverage.

1141 **Criminal Records:** Suppose a model f is trained on criminal records to predict individual crime
 1142 likelihood. Additionally, suppose the criminal record x_p of a person p is corrupted and publicly
 1143 available. Then, $f(x_p)$ predicts that person p is highly likely to commit crime. An adversarial law
 1144 enforcement agency, or even a prospective employer, may ignore or be unaware of warnings about
 1145 the data being corrupted, rendering a dangerous scenario for person p .⁵ To make this clear, provide a
 1146 figure similar to that of Fig. 1 at Fig. F.4.

1147 **Mortgage Loans:** Suppose a model f is trained on various items relevant to whether one receives
 1148 a mortgage loan or not, like bank statements and past rent payments. Person p has corrupted rent
 1149 payment history x_p . Then, the bank runs model f and obtains $f(x_p)$, which confidently says that x_p
 1150 is undeserving of a loan.

1151 **Car Insurance:** Suppose a model f is trained on one’s history of car accidents. Person p has
 1152 corrupted car accident history x_p . Then, when applying for car insurance, the provider runs model f
 1153 and obtains $f(x_p)$, which confidently says that x_p is undeserving of a loan.⁶

1154 We provide an additional example in the generative setting as well:

1155 **News Articles:** Consider a text-to-image generative model trained on a large dataset, including
 1156 web data, which has web articles and associated images. A popular news site publishes an article
 1157 about a businessperson, but mistakenly uses a picture of an unrelated individual p , x_p , as the header
 1158 image. This creates a strong, albeit false, association between this person’s likeness and the (perhaps
 1159 negative) content of the article. When prompted with a string similar to the headline of the news
 1160 article, the model generates an image (or a similar image) of person p , algorithmically cementing a
 1161 false narrative about person p .
 1162

1163 F LIMITATIONS AND FUTURE DIRECTIONS

1164 Notably, our presented method only applies to classification. While we provide an example of TTP in
 1165 the generative setting in Appx. E, we leave extending TTP to e.g. diffusion models (Song et al., 2021)
 1166 to future work. This extension is nontrivial, since there are uncountably infinitely many instances to
 1167 induce uniformity over. Furthermore, following the newspaper example given in Appx. E, the “forget
 1168 set” is often a *concept* or *association*, not a discrete set of instances. This holds similarly for
 1169 sequence-to-sequence generation (Vaswani et al., 2017). Furthermore, even in the discriminative
 1170 setting, we focus our method on image classification. Extending our methods to the text setting,
 1171 which is nontrivial due to discrete inputs, remains as future work.
 1172

1173 From an algorithmic perspective, in Alg. 1, we use linear scalarization to design our objective (Hwang
 1174 and Masud, 2012). One can instead design an objective using ε -constraints (Miettinen, 1999), which
 1175 can then be solved by an augmented Lagrangian method (Nocedal and Wright, 1999). Furthermore,
 1176 our certified algorithms, especially Alg. 3, are very sensitive to hyperparameters. Specifically, λ and
 1177 J must be nearly precisely $|\lambda_{\min}|$ and $|\lambda_{\max} + \lambda|$. Otherwise, the estimator diverges aggressively.
 1178 Furthermore, other Hessian estimators (Mehta et al., 2022) do not admit bounds which are necessary
 1179 for our certified Gaussian mechanism. Finally, taking λ large to ensure that the Hessian is p.d. and
 1180 hence invertible, i.e. local convex approximation, is problematic in our multiobjective setting since it
 1181 results in a step off the Pareto frontier. These limitations can be overcome by adopting recent Hessian-
 1182 free approaches from certified unlearning (Qiao et al., 2025; Mu and Klabjan, 2025; Koloskova et al.,
 1183 2025). However, adopting these approaches is nontrivial, since they rely on assumptions about the
 1184

1185 ⁵Recently, ML model providers have been involved in privacy cases involving criminal records (of Pri-
 1186 vacy Professionals, 2020), making this threat pertinent.

1187 ⁶Note that recent, the Department of Motor Vehicles in America has been selling driving records, making
 this threat pertinent (White, 2020).

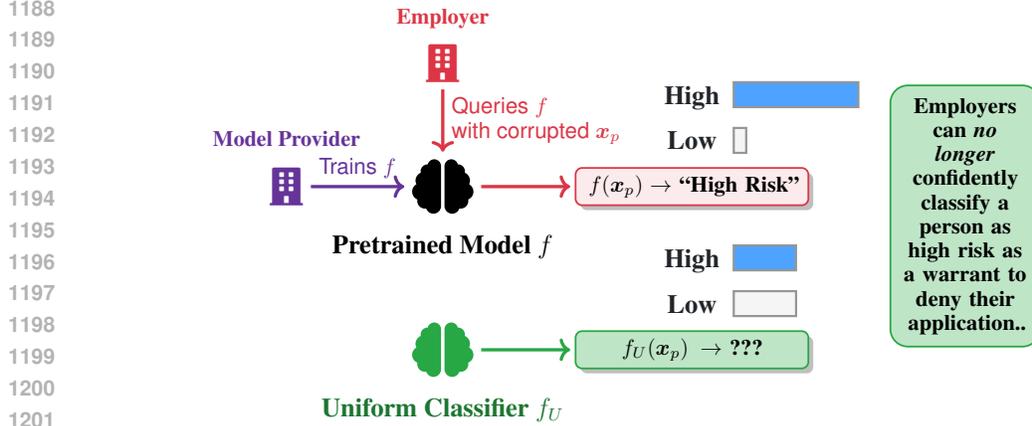


Figure F.4: An adversary, like an employer (🏢), can query a pretrained model f (🧠) and use its outputs to make harmful decisions. However, after running our algorithm, the new model f_U (🌱) provides maximal uncertainty, protecting against such an adversary. This is a duplicate of Fig. 1, to make clear how TTP extends to other settings.

retrain-from-scratch (target) model which do not hold for the Pareto model, e.g. the structure of its gradient updates. As such, we leave this to future work.

Finally, while we proved an upper bound on the retain accuracy drop for our Pareto learner without assuming convexity in Sec. 5, the general theoretical qualities of TTP still remain to be studied. In particular, one could prove lower bounds on the generalization error or the sample complexity of TTP. This would allow us to design more principled TTP algorithms which e.g. achieve the minimax rate. However, obtaining such lower bounds is nontrivial. For example, current approaches which prove lower bounds for deletion capacity (sample complexity) in the related setting of certified unlearning (Allouah et al., 2025) rely heavily on the fact that (ϵ, δ) -unlearning aims to mimic $\mathcal{A}(\mathcal{D}_r)$, whereas (ϵ, δ) -TTP aims to mimic $\mathcal{M}_\theta(\mathcal{D})$. Furthermore, these works assume strongly convex losses, making them inapplicable to the deep learning setting. We thus leave this to future work.

G DESIGNING CERTIFIED ALGORITHMS

In what follows, we design $(\epsilon, \delta, \theta)$ -certified Pareto learners. A symbol table can be found at Appx. N.

In our setting, the original model is obtained using ERM over some loss function \mathcal{L}_A , some dataset \mathcal{D} , and some parameter space \mathcal{W} . Furthermore, we consider the common scenario where the cumulative loss \mathcal{L}_A over the dataset is a finite sum of individual losses ℓ_A . Thus, we denote the pretrained model as:

$$\mathbf{w}^* = \mathcal{A}(\mathcal{D}) := \arg \min_{\mathbf{w} \in \mathcal{W}} \mathcal{L}_A(\mathbf{w}, \mathcal{D}) = \arg \min_{\mathbf{w} \in \mathcal{W}} \sum_{i=1}^{|\mathcal{D}|} \ell_A(\mathbf{w}, \mathcal{D}^{(i)}). \quad (\text{G.13})$$

By Prop. 4.2, we can similarly obtain a uniform learner through ERM with respect to some loss function \mathcal{L}_K . Furthermore, in our setting, we have the forget set \mathcal{D}_f and retain set $\mathcal{D}_r = \mathcal{D} \setminus \mathcal{D}_f$. Thus, the uniform learner over the forget set can be characterized as:

$$\mathcal{K}(\mathcal{D}_f) := \arg \min_{\mathbf{w} \in \mathcal{W}} \mathcal{L}_K(\mathbf{w}, \mathcal{D}) = \sum_{i=1}^{|\mathcal{D}_f|} \ell_K(\mathbf{w}, \mathcal{D}_f^{(i)}). \quad (\text{G.14})$$

Let $\theta \in (0, 1)$ be a tradeoff parameter between uniformity over the forget set and utility over the retain set. This yields a concrete characterization of \mathcal{M}_θ as:

1242

1243

1244

$$\tilde{\mathbf{w}}^* = \mathcal{M}_\theta(D) := \arg \min_{\mathbf{w} \in \mathcal{W}} \theta \mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r), \quad (\text{G.15})$$

1245

1246

1247

1248

$$= \arg \min_{\mathbf{w} \in \mathcal{W}} \theta \sum_{i=1}^{|\mathcal{D}_f|} \ell_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f^{(i)}) + (1 - \theta) \sum_{i=1}^{|\mathcal{D}_r|} \ell_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r^{(i)}). \quad (\text{G.16})$$

1249

as in Eq. (6).

1250

1251

1252

To design an algorithm which takes in \mathcal{D} , \mathcal{D}_r , and \mathbf{w}^* and outputs a parameter which satisfies Def. 4.6, we follow the methodology of certified unlearning Zhang et al. (2024), which seeks to satisfy Def. 4.5.

1253

First, we simplify the problem of deriving a model that satisfies Def. 4.6:

1254

1255

1256

1257

Theorem G.1. (Certification Guarantee) Let $\tilde{\mathbf{w}} := \mathcal{F}(\mathcal{D}, \mathcal{D}_f, \mathbf{w}^*)$ be an approximation to $\tilde{\mathbf{w}}^*$. Suppose $\|\tilde{\mathbf{w}} - \tilde{\mathbf{w}}^*\|_2 \leq \Delta$. Then, $\mathcal{U}(\mathcal{D}, \mathcal{D}_f, \mathcal{A}(D)) = \mathbf{w}^- = \tilde{\mathbf{w}} + \mathbf{Y}$ is a $(\epsilon, \delta, \theta)$ certified uniformity algorithm, where $\mathbf{Y} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ and $\sigma \geq \frac{\Delta}{\epsilon} \sqrt{2 \ln(1.25/\delta)}$.

1258

Proof: See Appx. H.4.

1259

1260

1261

Thus, it then suffices to find an approximation of $\tilde{\mathbf{w}}^*$, i.e. a form for $\mathcal{F}(\mathcal{D}, \mathcal{D}_f, \mathbf{w}^*)$ and its associated Δ . To do so, we consider the two assumptions Asm. 5.3 and Asm. 5.4.

1262

1263

1264

1265

For any $\mathbf{w} \in \mathcal{W}$, denote $\nabla_{\mathbf{w}, \mathcal{K}, \mathcal{A}} := \nabla_{\mathbf{w}}(\theta \mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r))$, the gradient of the objective of \mathcal{M}_θ with respect to \mathbf{w} , and $\mathbf{H}_{\mathbf{w}, \mathcal{K}, \mathcal{A}} := \nabla_{\mathbf{w}}^2(\theta \mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r))$, the Hessian of the objective of \mathcal{M}_θ with respect to \mathbf{w} . We thus have $\nabla_{\mathbf{w}, \mathcal{K}, \mathcal{A}} = \theta \nabla_{\mathbf{w}, \mathcal{K}} + (1 - \theta) \nabla_{\mathbf{w}, \mathcal{A}}$, and similarly for the Hessian.

1266

1267

Next, letting $g(\mathbf{w}) := \nabla_{\mathbf{w}, \mathcal{K}, \mathcal{A}}$, by Taylor's theorem, expanding $g(\tilde{\mathbf{w}}^*)$ around \mathbf{w}^* , we have that:

1268

1269

1270

$$g(\tilde{\mathbf{w}}^*) \approx g(\mathbf{w}^*) + Dg|_{\mathbf{w}^*}(\tilde{\mathbf{w}}^* - \mathbf{w}^*). \quad (\text{G.17})$$

1271

1272

1273

Note that $g(\tilde{\mathbf{w}}^*) = 0$, since $\tilde{\mathbf{w}}^*$ is the minimizer of the objective in \mathcal{M}_θ . Isolating $\tilde{\mathbf{w}}^*$ and using the definition of g , we then have that:

1274

1275

1276

$$\tilde{\mathbf{w}}^* \approx \mathbf{w}^* - \mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}^{-1} \nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}. \quad (\text{G.18})$$

1277

Thus, we let $\tilde{\mathbf{w}} = \mathbf{w}^* - \mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}^{-1} \nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}$. This yields the following general form of Δ :

1278

1279

1280

1281

1282

1283

1284

Proposition G.2. Suppose Asm. 5.3 and Asm. 5.4 hold. Suppose $\tilde{\mathbf{w}} = \mathbf{w}^* - \mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}^{-1} \nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}$. Then,

$$\|\tilde{\mathbf{w}}^* - \tilde{\mathbf{w}}\|_2 \leq \frac{\theta F_{\mathcal{K}} + (1 - \theta) F_{\mathcal{A}}}{2} \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}^{-1}\|_2 \|\mathbf{w}^* - \tilde{\mathbf{w}}^*\|_2^2. \quad (\text{G.19})$$

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

We then use local convex approximation (Nocedal and Wright, 1999) to bound $\|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}^{-1}\|_2$. To that end, we let the objective of \mathcal{M}_θ have a regularization term $\frac{\lambda}{2} \|\mathbf{w}\|_2^2$, yielding the inverse Hessian $\|(\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1}\|_2$; thus, in Prop. G.2, the norm of the inverse Hessian is replaced by $\|(\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1}\|_2$. It then suffices to bound this term.

Additionally, note that since the objective of \mathcal{M}_θ is nonconvex, the Hessian may not be invertible, i.e. $\lambda_{\min}(\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}) < 0$. However, $\lambda_{\min}(\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I}) = \lambda_{\min}(\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}) + \lambda$. Thus, for λ sufficiently large, we can make $\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I}$ positive definite and hence invertible, resolving this issue. In particular, we can take $\lambda > \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}\|_2$.

Furthermore, we let $\|\mathbf{w}\|_2 \leq C$ in \mathcal{M}_θ and \mathcal{A} , i.e. $\mathcal{M}_\theta = \arg \min_{\|\mathbf{w}\|_2 \leq C, \mathbf{w} \in \mathcal{W}} \theta \mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_A) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$ and $\mathcal{A}(\mathcal{D}) = \arg \min_{\mathbf{w} \in \mathcal{W}, \|\mathbf{w}\|_2 \leq C} \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D})$. Note that, as mentioned in (Zhang et al., 2024), unlearning methods implicitly assume this.

Together, these two methods yield a tractable form of Δ :

Proposition G.3. *Suppose Asm. 5.3 and Asm. 5.4 hold. Suppose $\tilde{\mathbf{w}} = \mathbf{w}^* - (\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} \nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}$ where $\|\mathbf{w}^*\|_2, \|\tilde{\mathbf{w}}^*\|_2 \leq C$. Let $\lambda_{\min} := \lambda_{\min}(\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}})$. Suppose $\lambda > \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}\|_2$. Then,*

$$\|\tilde{\mathbf{w}}^* - \tilde{\mathbf{w}}\|_2 \leq \frac{2C((\theta M_K + (1 - \theta)M_A)C + \lambda)}{\lambda + \lambda_{\min}}. \quad (\text{G.20})$$

Proof: See Appx. H.6.

While Prop. G.3 does yield a form of \mathcal{F} and Δ , the computation of $\tilde{\mathbf{w}}$ requires obtaining the exact inverse Hessian, which has runtime $\mathcal{O}(dz^2 + z^3)$, where z is the number of learnable parameters. Furthermore, computing the gradient product with the inverse Hessian is $\mathcal{O}(z^2)$. Finally, computing the gradient $\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}$ is $\mathcal{O}(|\mathcal{D}|z)$. Thus, the algorithm yielded by Prop. G.3 has a runtime complexity of $\mathcal{O}(dz^2 + z^3 + z^2 + |\mathcal{D}|z)$.

If we consider the additional assumption of convexity, we can take λ very small to ensure the Hessian is invertible, since we have $\lambda_{\min} = 0$. Thus, for convex models e.g. logistic regression with a mean-square uniform loss, this is tractable. This yields Alg. 2.

However, for nonconvex models e.g. large scale neural networks, this is computationally intractable. Thus, to provide better runtime, we derive an asymptotically unbiased estimator of the inverse Hessian. However, the estimator in Zhang et al. (2024) does not trivially extend to our case. In particular, we cannot glean Hessian samples using sampled i.i.d. data from the retain set, because the Hessian in our setting is defined over the forget set as well. Thus, we must derive an unbiased estimator while sampling Hessians from *both* the retain and forget set. As such, following the techniques of (Agarwal et al., 2016), we design an unbiased estimator as follows:

Theorem G.4. *Suppose we have n i.i.d. data samples (X_1, \dots, X_n) drawn from D_f and D_r , uniformly at random, with probabilities θ and $1 - \theta$ respectively. Then, suppose $\|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I}\|_2 \leq J$. For $t = 1, \dots, n$, if $X_t \sim D_f$ let $\mathbf{H}_{t, \lambda} = \mathbf{H}_{\mathbf{w}^*, \mathcal{K}, t} + \frac{\lambda \mathbf{I}}{2\theta}$ and if $X_t \sim D_r$ let $\mathbf{H}_{t, \lambda} = \mathbf{H}_{\mathbf{w}^*, \mathcal{A}, t} + \frac{\lambda \mathbf{I}}{2(1-\theta)}$. Suppose $\lambda > \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}\|_2$. Then, compute:*

$$\tilde{\mathbf{H}}_{t, \lambda}^{-1} = \mathbf{I} + (\mathbf{I} - \frac{\mathbf{H}_{t, \lambda}}{J})\tilde{\mathbf{H}}_{t-1, \lambda}^{-1}, \quad \tilde{\mathbf{H}}_{0, \lambda} = \mathbf{I}. \quad (\text{G.21})$$

Then, $\frac{\tilde{\mathbf{H}}_{n, \lambda}^{-1}}{J}$ is an asymptotically unbiased estimator for $(\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1}$

Proof: See Appx. H.7

One simple choice of J is $J = 2\lambda$, by Lemma H.1. However, we let J be free. The computation of the estimator in theorem G.4 has a runtime complexity of $\mathcal{O}(nz^2)$, a great speedup over the original $\mathcal{O}(dz^2 + z^3)$. Furthermore, with Hessian vector product (HVP) techniques (Pearlmutter, 1994), we obtain a space complexity of $\mathcal{O}(z)$ instead of $\mathcal{O}(z^2)$, since we do not have to compute the sample Hessians explicitly. Furthermore, computing $\tilde{\mathbf{H}}_{n, \lambda}^{-1} \nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}$ recursively reduces $\mathcal{O}(z^2)$ to $\mathcal{O}(nz)$.

Additionally, following Agarwal et al. (2016), we can average b unbiased estimators $\frac{\tilde{\mathbf{H}}_{t, \lambda}^{-1}}{J}$ as $\frac{1}{b} \sum_{i=1}^b \frac{\tilde{\mathbf{H}}_{t, \lambda}^{-1, (i)}}{J}$ to achieve better concentration. Altogether, we achieve a final runtime complexity of $\mathcal{O}(bnz^2 + bnz + |\mathcal{D}|z)$.

Furthermore, we relax the assumption that \mathbf{w}^* and $\tilde{\mathbf{w}}^*$ are the global minimizers of $\mathcal{L}_{\mathcal{A}}$ and $\theta \mathcal{L}_{\mathcal{A}} + (1 - \theta) \mathcal{L}_{\mathcal{K}}$. We do so because, in practice, it is possible that the data controller trained their model with early stopping, i.e. they did not reach the global minimizer. Altogether, this yields a final form of Δ as:

Theorem G.5. *Let $\tilde{\mathbf{w}}^*$ and \mathbf{w}^* not be empirical risk minimizers of their respective losses, but rather approximations thereof. Suppose Asm. 5.3 and Asm. 5.4 hold. Suppose $\|\mathbf{w}^*\|_2, \|\tilde{\mathbf{w}}^*\|_2 \leq C$. Let $\lambda_{\min} := \lambda_{\min}(\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}})$. Suppose $\lambda > \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}\|_2$. Let $\tilde{\mathbf{w}} = \mathbf{w}^* - \frac{\tilde{\mathbf{H}}_{t, \lambda}^{-1}}{J} \nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}$. Let*

Algorithm 3 $(\epsilon, \delta, \theta)$ -Certified Uniformity with Inverse Hessian Estimator

Require: Dataset \mathcal{D} ; forget set \mathcal{D}_f ; pretrained model $w^* = \mathcal{A}(\mathcal{D})$; privacy budgets ϵ and δ ; uniformity-utility tradeoff coefficient θ ; estimator concentration b ; sample size n ; local convex coefficient λ ; norm upper bound C ; cumulative Hessian upper bound H ; individual Hessian minimum eigenvalue upper bound ζ_{\min} ; gradient norm upper bound G ; bound looseness probability ρ .

```

1356  $P_{0,\lambda}^{(0)} \leftarrow \nabla_{w^*, \mathcal{K}, \mathcal{A}}$ 
1357 for  $j = 1, \dots, b$  do
1358   for  $t = 1, \dots, n$  do
1359     Sample  $X_t$  from  $D_f$  uniformly with probability  $\theta$  or,
1360     sample  $X_t$  from  $D_r$  uniformly with probability  $1 - \theta$ .
1361     if  $X_t \sim D_f$  then
1362        $H_{t,\lambda}^{(j)} \leftarrow \nabla_w^2 \mathcal{L}_{\mathcal{K}}(w^*, X_i) + \frac{\lambda I}{2\theta}$ .
1363     else if  $X_t \sim D_r$  then
1364        $H_{t,\lambda}^{(j)} \leftarrow \nabla_w^2 \mathcal{L}_{\mathcal{A}}(w^*, X_i) + \frac{\lambda I}{2(1-\theta)}$ .
1365     end if
1366      $P_{t,\lambda}^{(j)} = P_{0,\lambda}^{(0)} + (I - \frac{H_{t,\lambda}^{(j)}}{H})P_{t-1,\lambda}^{(j)}$ .
1367   end for
1368 end for
1369  $P_{n,\lambda} \leftarrow \frac{1}{b} \sum_{j=1}^b P_{n,\lambda}^{(j)}$ .
1370  $\tilde{w} \leftarrow w^* - \frac{P_{n,\lambda}}{H}$ .
1371 Compute  $\Delta$  as the bound in Eq. (G.23).
1372  $\sigma = \frac{\Delta}{\epsilon} \sqrt{2 \ln(1.25/\delta)}$ 
1373  $w^- \leftarrow \tilde{w} + Y$  where  $Y \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$ .
1374 return  $w^-$ .

```

b be the number of inverse Hessian estimators we average. Letting n be the number of steps taken during unbiased estimation of the inverse Hessian, require $n \geq 2 \frac{B}{\lambda + \lambda_{\min}} \ln(\frac{B}{\lambda + \lambda_{\min}} b)$ where $B = \max\{\frac{\theta P_{\mathcal{K}} + \lambda}{|\mathcal{D}_f|}, \frac{(1-\theta)P_{\mathcal{A}} + \lambda}{|\mathcal{D}_r|}\}$. Suppose $\|\nabla_{w^*, \mathcal{K}, \mathcal{A}}\|_2, \|\nabla_{\tilde{w}^*, \mathcal{K}, \mathcal{A}}\|_2 \leq G$, With probability larger than $1 - \rho$, we have that:

$$\|\tilde{w}^* - \tilde{w}\|_2 \leq \frac{2C((\theta F_{\mathcal{K}} + (1-\theta)F_{\mathcal{A}})C + \lambda) + G}{\lambda + \lambda_{\min}} \quad (\text{G.22})$$

$$+ (16 \frac{B}{\zeta_{\min}} \sqrt{\frac{\ln(\frac{d}{\rho})}{b}} + \frac{1}{16})(2C(\theta P_{\mathcal{K}} + (1-\theta)P_{\mathcal{A}}) + G). \quad (\text{G.23})$$

where $\zeta_{\min} \geq \min_i \lambda_{\min}(\nabla_w^2 \tilde{\ell}_{\mathcal{K}, \mathcal{A}}(w, \mathcal{D}^{(i)}))$.

Proof: See Appx. H.8.

Note that if we let w^* be an ERM in theorem G.5, we can use $\nabla_{w, \mathcal{K}, \mathcal{A}}$ and obtain the same result. Altogether, this yields Alg. 3.

H PROOFS

H.1 HELPFUL LEMMAS

Lemma H.1. Given Asm. 5.3, the gradients $\nabla_{w, \mathcal{K}}$ and $\nabla_{w, \mathcal{A}}$ exist and are Lipschitz with constants $P_{\mathcal{K}}$ and $P_{\mathcal{A}}$, respectively. Furthermore, given Asm. 5.4, the Hessians $\mathbf{H}_{w, \mathcal{K}}$ and $\mathbf{H}_{w, \mathcal{A}}$ exist and are Lipschitz with constants $F_{\mathcal{K}}$ and $F_{\mathcal{A}}$, respectively.

1404 *Proof.*

$$1405 \|\nabla_{\mathbf{w}_1, \mathcal{K}} - \nabla_{\mathbf{w}_2, \mathcal{K}}\|_2 = \left\| \sum_{i=1}^{|\mathcal{D}_f|} \nabla^2 \ell_{\mathcal{K}}^{(i)}(\mathbf{w}_1, \mathcal{D}_f^{(i)}) - \sum_{i=1}^{|\mathcal{D}_f|} \nabla^2 \ell_{\mathcal{K}}^{(i)}(\mathbf{w}_2, \mathcal{D}_f^{(i)}) \right\|_2 \quad (\text{H.24})$$

$$1409 \leq \sum_{i=1}^{|\mathcal{D}_f|} \|\nabla^2 \ell_{\mathcal{K}}^{(i)}(\mathbf{w}_1, \mathcal{D}_f^{(i)}) - \nabla^2 \ell_{\mathcal{K}}^{(i)}(\mathbf{w}_2, \mathcal{D}_f^{(i)})\|_2, \text{ triangle inequality} \quad (\text{H.25})$$

$$1412 \leq \sum_{i=1}^{|\mathcal{D}_f|} \frac{P_{\mathcal{K}}}{|\mathcal{D}_f|} \|\mathbf{w}_1 - \mathbf{w}_2\|_2, \text{ Asm. 5.3} \quad (\text{H.26})$$

$$1415 = P_{\mathcal{K}} \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad (\text{H.27})$$

1416 This follows similarly for $\nabla_{\mathbf{w}, \mathcal{A}}$, $\mathbf{H}_{\mathbf{w}, \mathcal{K}}$, and $\mathbf{H}_{\mathbf{w}, \mathcal{A}}$. \square

1418 **Lemma H.2.** *Given Asm. 5.3, for any dataset $\mathcal{D} \subset \mathcal{Z}^n$, $\mathcal{L}_{\mathcal{A}}$ satisfies:*

$$1421 |\mathcal{L}_{\mathcal{A}}(\mathbf{w}_1, \mathcal{D}) - \mathcal{L}_{\mathcal{A}}(\mathbf{w}_2, \mathcal{D})| \leq \frac{P_{\mathcal{K}}}{2} \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 + \|\nabla_{\mathbf{w}_2, \mathcal{A}}\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad (\text{H.28})$$

1424 *Proof.* By the fundamental theorem of calculus, we have the path integral:

$$1426 \int_0^1 \langle \nabla_{\mathbf{w}_2 + t(\mathbf{w}_1 - \mathbf{w}_2), \mathcal{A}}, \mathbf{w}_1 - \mathbf{w}_2 \rangle dt = \mathcal{L}_{\mathcal{A}}(\mathbf{w}_1, \mathcal{D}) - \mathcal{L}_{\mathcal{A}}(\mathbf{w}_2, \mathcal{D}) \quad (\text{H.29})$$

1429 We have that:

$$1432 \int_0^1 \langle \nabla_{\mathbf{w}_2 + t(\mathbf{w}_1 - \mathbf{w}_2), \mathcal{A}}, \mathbf{w}_1 - \mathbf{w}_2 \rangle dt = \int_0^1 \langle \nabla_{\mathbf{w}_2, \mathcal{A}} - \nabla_{\mathbf{w}_2, \mathcal{A}} + \nabla_{\mathbf{w}_2 + t(\mathbf{w}_1 - \mathbf{w}_2), \mathcal{A}}, \mathbf{w}_1 - \mathbf{w}_2 \rangle dt \quad (\text{H.30})$$

$$1435 = \int_0^1 \langle \nabla_{\mathbf{w}_2, \mathcal{A}}, \mathbf{w}_1 - \mathbf{w}_2 \rangle dt \quad (\text{H.31})$$

$$1438 + \int_0^1 \langle \nabla_{\mathbf{w}_2 + t(\mathbf{w}_1 - \mathbf{w}_2), \mathcal{A}} - \nabla_{\mathbf{w}_2, \mathcal{A}}, \mathbf{w}_1 - \mathbf{w}_2 \rangle dt \quad (\text{H.32})$$

1440 The first term can be bounded by Cauchy-Schwarz as:

$$1442 \int_0^1 \langle \nabla_{\mathbf{w}_2, \mathcal{A}}, \mathbf{w}_1 - \mathbf{w}_2 \rangle dt \leq \|\nabla_{\mathbf{w}_2, \mathcal{A}}\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad (\text{H.33})$$

1446 and similarly the second term can be bounded by Cauchy-Schwarz as:

$$1448 \int_0^1 \langle \nabla_{\mathbf{w}_2 + t(\mathbf{w}_1 - \mathbf{w}_2), \mathcal{A}} - \nabla_{\mathbf{w}_2, \mathcal{A}}, \mathbf{w}_1 - \mathbf{w}_2 \rangle dt \leq \int_0^1 \|\nabla_{\mathbf{w}_2 + t(\mathbf{w}_1 - \mathbf{w}_2), \mathcal{A}} - \nabla_{\mathbf{w}_2, \mathcal{A}}\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 dt \quad (\text{H.34})$$

$$1452 \leq \int_0^1 P_{\mathcal{K}} t \|\mathbf{w}_1 - \mathbf{w}_2\|_2, \text{ by Lemma H.1} \quad (\text{H.35})$$

$$1454 \leq \frac{P_{\mathcal{A}}}{2} \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad (\text{H.36})$$

1456 Incorporating these bounds into Eq. (H.32) and Eq. (H.29), upon applying the triangle inequality, yields:

$$|\mathcal{L}_A(\mathbf{w}_1, \mathcal{D}) - \mathcal{L}_A(\mathbf{w}_2, \mathcal{D})| \leq \frac{P_{\mathcal{K}}}{2} \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 + \|\nabla_{\mathbf{w}_2, \mathcal{A}}\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad (\text{H.37})$$

as desired. \square

Lemma H.3. *Given Asm. 5.4, the Hessians $\mathbf{H}_{\mathbf{w}, \mathcal{K}}$, $\mathbf{H}_{\mathbf{w}, \mathcal{A}}$, and $\mathbf{H}_{\mathbf{w}, \mathcal{K}, \mathcal{A}}$ are symmetric.*

Proof. By Lemma H.1, the Hessians $\mathbf{H}_{\mathbf{w}, \mathcal{K}}$ and $\mathbf{H}_{\mathbf{w}, \mathcal{A}}$ are continuous, and thus $\mathbf{H}_{\mathbf{w}, \mathcal{K}, \mathcal{A}}$ is continuous by linearity. Hence, all second-order partial derivatives contained in the Hessians are continuous, so by Schwartz's theorem all Hessians are symmetric. Importantly, for e.g. $\mathbf{H}_{\mathbf{w}, \mathcal{K}, \mathcal{A}}$, $\|\mathbf{H}_{\mathbf{w}, \mathcal{K}, \mathcal{A}}\|_2 = \max_i |\lambda_i(\mathbf{H}_{\mathbf{w}, \mathcal{K}, \mathcal{A}})|$, where λ_i denotes the i th eigenvalue. \square

Lemma H.4. *(Corollary of Theorem A.1 in (Dwork et al., 2014)) Let $X \sim \mathcal{N}(\lambda, \sigma^2 \mathbf{I})$ and $Y \sim \mathcal{N}(\lambda', \sigma^2 \mathbf{I})$. Suppose $\|\lambda - \lambda'\|_2 \leq \Delta$. Then for any $\delta > 0$, X and Y are (ε, δ) -indistinguishable if $\sigma \geq \frac{\Delta}{\varepsilon} \sqrt{2 \ln(1.25/\delta)}$.*

Lemma H.5. *Suppose we have n i.i.d. data samples (X_1, \dots, X_n) drawn from D_f and D_r with probabilities θ and $1 - \theta$ respectively. For $t = 1, \dots, n$, if $X_t \sim D_f$ let $\mathbf{H}_{t, \lambda} = \mathbf{H}_{\mathbf{w}^*, \mathcal{K}, t} + \frac{\lambda \mathbf{I}}{2\theta}$ and if $X_t \sim D_r$ let $\mathbf{H}_{t, \lambda} = \mathbf{H}_{\mathbf{w}^*, \mathcal{A}, t} + \frac{\lambda \mathbf{I}}{2(1-\theta)}$. Then, $\mathbb{E}[\mathbf{H}_{t, \lambda}] = \mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I}$ i.e. $\mathbf{H}_{t, \lambda}$ is an unbiased estimator of our Hessian of interest.*

Proof. At time t , we have sample X_t s.t. $X_t \sim D_r$ or $X_t \sim D_f$. Note that $\mathbb{E}_{X_t \sim D_f}[\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, t} + \frac{\lambda \mathbf{I}}{2\theta}] = \mathbf{H}_{\mathbf{w}^*, \mathcal{K}} + \frac{\lambda \mathbf{I}}{2\theta}$, and likewise $\mathbb{E}_{X_t \sim D_r}[\mathbf{H}_{\mathbf{w}^*, \mathcal{A}, t} + \frac{\lambda \mathbf{I}}{2(1-\theta)}] = \mathbf{H}_{\mathbf{w}^*, \mathcal{A}} + \frac{\lambda \mathbf{I}}{2(1-\theta)}$.

By the law of iterated expectation, we have that:

$$\begin{aligned} \mathbb{E}[\mathbf{H}_{t, \lambda}] &= \mathbb{E}[\mathbf{H}_{t, \lambda} | X_t \sim D_f] \Pr(X_t \sim D_f) + \mathbb{E}[\mathbf{H}_{t, \lambda} | X_t \sim D_r] \Pr(X_t \sim D_r) \\ &= \theta \mathbb{E}_{X_t \sim D_f}[\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, t} + \frac{\lambda \mathbf{I}}{2\theta}] + (1 - \theta) \mathbb{E}_{X_t \sim D_r}[\mathbf{H}_{\mathbf{w}^*, \mathcal{A}, t} + \frac{\lambda \mathbf{I}}{2(1-\theta)}] \\ &= \theta(\mathbf{H}_{\mathbf{w}^*, \mathcal{K}} + \frac{\lambda \mathbf{I}}{2\theta}) + (1 - \theta)(\mathbf{H}_{\mathbf{w}^*, \mathcal{A}} + \frac{\lambda \mathbf{I}}{2(1-\theta)}) \\ &= \theta \mathbf{H}_{\mathbf{w}^*, \mathcal{K}} + (1 - \theta) \mathbf{H}_{\mathbf{w}^*, \mathcal{A}} + \lambda \mathbf{I} \\ &= \mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I} \end{aligned}$$

as desired. \square

Lemma H.6. *Suppose Assumptions 5.3 and 5.4 hold. Let local condition number $\hat{\kappa}_l$ and maximum local condition number $\hat{\kappa}_l^{\max}$ correspond to the definitions of Agarwal et al. (2016) with respect to the Hessian of the loss of \mathcal{M}_θ after local convex approximation. Then, $\hat{\kappa}_l \leq \frac{B}{\lambda + \lambda_{\min}}$ and $\hat{\kappa}_l^{\max} \leq \frac{B}{\zeta_{\min}}$ where $B = \max\{\frac{\theta P_{\mathcal{K}} + \lambda}{|\mathcal{D}_f|}, \frac{(1-\theta)P_{\mathcal{A}} + \lambda}{|\mathcal{D}_r|}\}$ and where $\zeta_{\min} \geq \min_i \lambda_{\min}(\nabla_{\mathbf{w}}^2 \tilde{\ell}_{\mathcal{K}, \mathcal{A}}^{(i)}(\mathbf{w}, \mathcal{D}^{(i)}))$.*

Proof. By Eq. (6) and our local convex approximation technique, we have that:

$$M_\theta(D) = \arg \min_{\mathbf{w} \in \mathcal{W}} \theta \sum_{i=1}^{|\mathcal{D}_f|} \ell_{\mathcal{K}}^{(i)}(\mathbf{w}, \mathcal{D}_f^{(i)}) + (1 - \theta) \sum_{i=1}^{|\mathcal{D}_r|} \ell_{\mathcal{A}}^{(i)}(\mathbf{w}, \mathcal{D}_r^{(i)}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (\text{H.38})$$

$$= \arg \min_{\mathbf{w} \in \mathcal{W}} \sum_{i=1}^{|\mathcal{D}_f|} (\theta \ell_{\mathcal{K}}^{(i)}(\mathbf{w}, \mathcal{D}_f^{(i)}) + \frac{\lambda}{2|\mathcal{D}_f|} \|\mathbf{w}\|_2^2) + \sum_{i=1}^{|\mathcal{D}_r|} (\theta \ell_{\mathcal{A}}^{(i)}(\mathbf{w}, \mathcal{D}_r^{(i)}) + \frac{\lambda}{2|\mathcal{D}_r|} \|\mathbf{w}\|_2^2) \quad (\text{H.39})$$

$$= \arg \min_{\mathbf{w} \in \mathcal{W}} \sum_{i=1}^{|\mathcal{D}|} \tilde{\ell}_{\mathcal{K}, \mathcal{A}}^{(i)}(\mathbf{w}, \mathcal{D}^{(i)}) \quad (\text{H.40})$$

1512 where

$$1513 \tilde{\ell}_{\mathcal{K},\mathcal{A}}^{(i)}(\mathbf{w}, \mathcal{D}^{(i)}) = \begin{cases} \theta \ell_{\mathcal{K}}^{(i)}(\mathbf{w}, \mathcal{D}_f^{(i)}) + \frac{\lambda}{2|\mathcal{D}_f|} \|\mathbf{w}\|_2^2, & 1 \leq i \leq |\mathcal{D}_f| \\ (1-\theta) \ell_{\mathcal{K}}^{(i-|\mathcal{D}_f|)}(\mathbf{w}, \mathcal{D}_r^{(i-|\mathcal{D}_f|)}) + \frac{\lambda}{2|\mathcal{D}_r|} \|\mathbf{w}\|_2^2, & |\mathcal{D}_f| + 1 \leq i \leq |\mathcal{D}| \end{cases} \quad (\text{H.41})$$

1518 By the definitions provided in Agarwal et al. (2016), we have:

$$1519 \hat{\kappa}_l = \max_{\mathbf{w} \in \mathcal{W}} \frac{\max_i \lambda_{\max}(\nabla_{\mathbf{w}}^2 \tilde{\ell}_{\mathcal{K},\mathcal{A}}(\mathbf{w}, \mathcal{D}^{(i)}))}{\lambda_{\min}(\mathbf{H}_{\mathbf{w},\mathcal{K},\mathcal{A}} + \lambda \mathbf{I})} \quad (\text{H.42})$$

1524 and

$$1525 \hat{\kappa}_l^{\max} = \max_{\mathbf{w} \in \mathcal{W}} \frac{\max_i \lambda_{\max}(\nabla_{\mathbf{w}}^2 \tilde{\ell}_{\mathcal{K},\mathcal{A}}(\mathbf{w}, \mathcal{D}^{(i)}))}{\min_i \lambda_{\min}(\nabla_{\mathbf{w}}^2 \tilde{\ell}_{\mathcal{K},\mathcal{A}}(\mathbf{w}, \mathcal{D}^{(i)}))} \quad (\text{H.43})$$

1529 We then have that, for any i ,

$$1530 \lambda_{\max}(\nabla_{\mathbf{w}}^2 \tilde{\ell}_{\mathcal{K},\mathcal{A}}^{(i)}) \leq \|\nabla_{\mathbf{w}}^2 \tilde{\ell}_{\mathcal{K},\mathcal{A}}^{(i)}\|_2, \text{ by Lemma H.3} \quad (\text{H.44})$$

$$1531 = \max\{\|\theta \nabla_{\mathbf{w}}^2 \ell_{\mathcal{K}}^{(i)} + \frac{\lambda \mathbf{I}}{|\mathcal{D}_f|}\|_2, \|(1-\theta) \nabla_{\mathbf{w}}^2 \ell_{\mathcal{A}}^{(i)} + \frac{\lambda \mathbf{I}}{|\mathcal{D}_r|}\|_2\} \quad (\text{H.45})$$

$$1532 \quad (\text{H.46})$$

1537 Furthermore, by Asm. 5.3 and the triangle inequality:

$$1538 \|\theta \nabla_{\mathbf{w}}^2 \ell_{\mathcal{K}}^{(i)} + \frac{\lambda \mathbf{I}}{|\mathcal{D}_f|}\|_2 \leq \frac{\theta P_{\mathcal{K}} + \lambda}{|\mathcal{D}_f|} \quad (\text{H.47})$$

1543 and

$$1544 \|(1-\theta) \nabla_{\mathbf{w}}^2 \ell_{\mathcal{A}}^{(i)} + \frac{\lambda \mathbf{I}}{|\mathcal{D}_r|}\|_2 \leq \frac{(1-\theta) P_{\mathcal{A}} + \lambda}{|\mathcal{D}_r|} \quad (\text{H.48})$$

1547 Taking max over all i , we obtain that

$$1548 \max_i \lambda_{\max}(\nabla_{\mathbf{w}}^2 \tilde{\ell}_{\mathcal{K},\mathcal{A}}(\mathbf{w}, \mathcal{D}^{(i)})) \leq \max\left\{\frac{\theta P_{\mathcal{K}} + \lambda}{|\mathcal{D}_f|}, \frac{(1-\theta) P_{\mathcal{A}} + \lambda}{|\mathcal{D}_r|}\right\} \quad (\text{H.49})$$

1551 which we denote by B .

1552 Then, we obtain that $\hat{\kappa}_l \leq \frac{B}{\lambda + \lambda_{\min}}$ and $\hat{\kappa}_l^{\max} \leq \frac{B}{\zeta_{\min}}$ as desired, since

1555 \square

1556 **Lemma H.7.** (Lemma 3.6 adapted from Agarwal et al. (2016)) Suppose Asm. 5.3 and Asm. 5.4 hold.

1557 Consider the estimator $\frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H}$ in theorem G.4. Let b be the number of inverse Hessian estimators we
1558 obtain. Suppose $n \geq 2 \frac{B}{\lambda + \lambda_{\min}} \ln(\frac{B}{\lambda + \lambda_{\min}} b)$, where $B = \max\{\frac{\theta P_{\mathcal{K}} + \lambda}{|\mathcal{D}_f|}, \frac{(1-\theta) P_{\mathcal{A}} + \lambda}{|\mathcal{D}_r|}\}$. Then, we have
1559 that:

$$1560 \Pr[\|\mathbf{H}_{\mathbf{w}^*,\mathcal{K},\mathcal{A}} + \lambda \mathbf{I}\|^{-1} - \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H}\|_2 \leq 16 \frac{B}{\zeta_{\min}} \sqrt{\frac{\ln(\frac{d}{\rho})}{b}} + \frac{1}{16}] \geq 1 - \rho \quad (\text{H.50})$$

1561 where $\zeta_{\min} \geq \min_i \lambda_{\min}(\nabla_{\mathbf{w}}^2 \tilde{\ell}_{\mathcal{K},\mathcal{A}}^{(i)}(\mathbf{w}, \mathcal{D}^{(i)}))$.

Proof. Note that $b = S_1$ in our setting. In our setting, following the subsequent steps of the proof in Agarwal et al. (2016) after plugging in the bounds in Lemma H.6 in place of $\hat{\kappa}_l, \hat{\kappa}_l^{\max}$, noting that we choose $n = S_2 \geq 2 \frac{B}{\lambda + \lambda_{\min}} \ln(\frac{B}{\lambda + \lambda_{\min}} b)$, we obtain the exact same result for the Neumann series bound of $\frac{1}{16}$. Using the fact that $\frac{B}{\zeta_{\min}}$ is an upper bound on $\hat{\kappa}_l^{\max}$ by Lemma H.6, the rest of the proof follows similarly. \square

Lemma H.8. (Proposition 2.1 in Dwork et al. (2014)) Let $\mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \rightarrow R$ be a randomized algorithm that is (ε, δ) -differentially private. Let $f : R \rightarrow R'$ be an arbitrary mapping. Then, $f \circ \mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \rightarrow R'$ is (ε, δ) -differentially private.

Note that, in the proof of Lemma H.8, one proves this fact for deterministic mappings, so this holds for both randomized and deterministic f .

Lemma H.9. Consider the mapping $\mathcal{J} : \mathcal{W} \rightarrow R$, and suppose $\mathcal{G} : \mathcal{Z}^n \times \mathcal{Z}^n \times \mathcal{W} \rightarrow \mathcal{W}$ satisfies Def. 4.6. Then, $\forall C \subset R$:

$$\Pr(\mathcal{J}(\mathcal{G}(\mathcal{D}, \mathcal{D}_f, \mathcal{A}(\mathcal{D}))) \in C) \leq e^\varepsilon \Pr(\mathcal{J}(\mathcal{M}_\theta(\mathcal{D})) \in C) + \delta \quad (\text{H.51})$$

$$\Pr(\mathcal{J}(\mathcal{M}_\theta(\mathcal{D})) \in C) \leq e^\varepsilon \Pr(\mathcal{J}(\mathcal{G}(\mathcal{D}, \mathcal{D}_f, \mathcal{A}(\mathcal{D}))) \in C) + \delta \quad (\text{H.52})$$

Proof. Immediate from Lemma H.8. \square

H.2 PROOF OF PROPOSITION 4.2

Proof. Fix a dataset $\mathcal{D} \subset \mathcal{Z}^n$.

Suppose we have an K -layer function $f_{\mathbf{w}} : \mathbb{R}^d \rightarrow \mathbb{R}^o$ parameterized by $\mathbf{w} \in \mathcal{W}$ of the form $f(\mathbf{x}) = L_1 \circ \dots \circ L_{K-1} \circ L_K$ where $L_{K-1}(\mathbf{x}) = \mathbf{W}_{K-1}^T \mathbf{x} + \mathbf{b}_{K-1}$ and $L_K(\mathbf{x}) = \text{softmax}(\mathbf{x})$, i.e. $L_{K-1}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{j=1}^{|\mathcal{Y}|} e^{x_j}}$. Thus, $f_{\mathbf{w}} \in \mathcal{H}_{\mathcal{W}}$. Then, let $\mathbf{W}_{K-1} = \mathbf{0}$ and $\mathbf{b}_{K-1} = \mathbf{0}$.

Fix $\mathbf{z} \in \mathcal{D}$. This yields, for $j = 1, \dots, |\mathcal{Y}|$, $f(\mathbf{z})_j = \frac{e^0}{\sum_{j=1}^{|\mathcal{Y}|} e^0} = \frac{e^0}{|\mathcal{Y}|e^0} = \frac{1}{|\mathcal{Y}|}$. Hence, since \mathbf{z} was

arbitrary, $f_{\mathbf{w}}(\mathbf{z}) = \underbrace{\left(\frac{1}{|\mathcal{Y}|}, \dots, \frac{1}{|\mathcal{Y}|} \right)}_{|\mathcal{Y}| \text{ times}} \forall \mathbf{z} \in \mathcal{D}$. Since \mathcal{D} was arbitrary, by definition of a uniform learner

over \mathcal{D} , $f_{\mathcal{K}(\mathcal{D})} \in \mathcal{H}_{\mathcal{W}} \forall \mathcal{D} \subset \mathcal{Z}^n$ as desired. \square

H.3 PROOF OF PROPOSITION 4.3

We use the following definition of global Pareto optimality:

Definition H.10. (Chapter 1 of Pardalos et al. (2017)) Suppose we have a multiobjective optimization problem $\min \mathbf{f}(\mathbf{x})$ s.t. $\mathbf{x} \in A$, where $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x}))$. $\mathbf{x}^* \in A$ with $\mathbf{f}(\mathbf{x}^*)$ is called globally Pareto optimal if and only if there exists no $\mathbf{x} \in A$ such that $f_i(\mathbf{x}) \leq f_i(\mathbf{x}^*)$ for all $i = 1, 2, \dots, m$ and $f_j(\mathbf{x}) < f_j(\mathbf{x}^*)$ for at least one $j \in \{1, \dots, m\}$.

We can then prove the statement:

Proof. Let $\theta \in (0, 1)$. Fix $\mathcal{D} \subset \mathcal{Z}^n$, $\mathcal{D}_f \subset \mathcal{D}$, and $\mathcal{D}_r = \mathcal{D} \setminus \mathcal{D}_f$.

Suppose, for the sake of contradiction, that $\tilde{\mathbf{w}}^* = \mathcal{M}_\theta(\mathcal{D}) = \text{argmin}_{\mathbf{w}} \theta \mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r)$, a global minimizer, is not globally Pareto optimal with respect to $\mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f)$ and $\mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r)$. Then, exists \mathbf{w}' s.t. $\mathcal{L}_{\mathcal{K}}(\mathbf{w}', \mathcal{D}_f) \leq \mathcal{L}_{\mathcal{K}}(\tilde{\mathbf{w}}^*, \mathcal{D}_f)$ and $\mathcal{L}_{\mathcal{A}}(\mathbf{w}', \mathcal{D}_r) \leq \mathcal{L}_{\mathcal{A}}(\tilde{\mathbf{w}}^*, \mathcal{D}_r)$, with at least one of these inequalities being strict.

Then, since $\theta \in (0, 1)$ and $(1 - \theta) \in (0, 1)$, we have that $\theta \mathcal{L}_{\mathcal{K}}(\mathbf{w}', \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathbf{w}', \mathcal{D}_r) < \theta \mathcal{L}_{\mathcal{K}}(\tilde{\mathbf{w}}^*, \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\tilde{\mathbf{w}}^*, \mathcal{D}_r)$, contradicting optimality of $\tilde{\mathbf{w}}^*$. As such, $\mathcal{M}_\theta(\mathcal{D})$ is globally Pareto optimal respect to $\mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f)$ and $\mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r)$ as desired.

This holds similarly for a local minimizer $\tilde{\mathbf{w}}^*$, where Pareto optimality similarly holds only locally in a neighborhood around the minima. \square

H.4 PROOF OF THEOREM G.1

Proof. The proof follows similarly to Lemma 10 in Sekhari et al. (2021); for completeness, we adapt their proof to our setting.

Let $\mathbf{w}^* := A(D)$, $\mathbf{w}^- := \mathcal{G}(\mathcal{D}, \mathcal{D}_f, \mathbf{w}^*)$, $\tilde{\mathbf{w}} := \mathcal{F}(\mathcal{D}, \mathcal{D}_f, \mathbf{w}^*)$. Departing from the notation of the theorem for clarity, let $\hat{\mathbf{w}}^* := \mathcal{M}_\theta(\mathcal{D})$, $\hat{\mathbf{w}}^- := \mathcal{G}(\mathcal{D}, \emptyset, \hat{\mathbf{w}}^*)$, $\hat{\tilde{\mathbf{w}}} := \mathcal{F}(\mathcal{D}, \emptyset, \hat{\mathbf{w}}^*)$.

Note that $\hat{\tilde{\mathbf{w}}} = \hat{\mathbf{w}}^*$. We then have that $\|\tilde{\mathbf{w}} - \hat{\tilde{\mathbf{w}}}\|_2 = \|\tilde{\mathbf{w}} - \hat{\mathbf{w}}^*\|_2 \leq \Delta$, by definition of Δ .

By definition of \mathcal{G} , we have that $\mathbf{w}^- = \tilde{\mathbf{w}} + Y$ and $\hat{\mathbf{w}}^- = \hat{\tilde{\mathbf{w}}} + Y$, where $Y \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ s.t. $\sigma \geq \frac{\Delta}{\varepsilon} \sqrt{2 \ln(1.25/\delta)}$.

As such, $\mathbf{w}^- = \mathcal{N}(\tilde{\mathbf{w}}, \sigma^2 \mathbf{I})$ and $\hat{\mathbf{w}}^- \sim \mathcal{N}(\hat{\tilde{\mathbf{w}}}, \sigma^2 \mathbf{I})$.

Thus, by Lemma H.4, \mathbf{w}^- , $\hat{\mathbf{w}}^-$ are (ε, δ) -indistinguishable. In particular, since $\hat{\mathbf{w}}^- = \hat{\mathbf{w}}^*$ by construction, $\mathcal{G}(\mathcal{D}, \mathcal{D}_f, \mathcal{A}(\mathcal{D}))$ and $\mathcal{M}_\theta(\mathcal{D})$ are (ε, δ) -indistinguishable, as desired. \square

H.5 PROOF OF PROPOSITION G.2

Proof. By the same token as Lemma 3.3 in (Zhang et al., 2024), we have that:

$$\|\tilde{\mathbf{w}} - \tilde{\mathbf{w}}^*\|_2 \leq \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}^{-1}\|_2 \int_0^1 \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \mathbf{H}_{\mathbf{w}^* + t(\tilde{\mathbf{w}}^* - \mathbf{w}^*), \mathcal{K}, \mathcal{A}}\|_2 \|\mathbf{w}^* - \tilde{\mathbf{w}}^*\|_2 dt \quad (\text{H.53})$$

Let $\mathbf{w}' = \mathbf{w}^* + t(\tilde{\mathbf{w}}^* - \mathbf{w}^*)$. We have that $\|\mathbf{w}^* - \mathbf{w}'\|_2 = \|\mathbf{w}^* - \mathbf{w}^* + t(\tilde{\mathbf{w}}^* - \mathbf{w}^*)\|_2 = t\|\mathbf{w}^* - \tilde{\mathbf{w}}^*\|_2$.

Furthermore, by linearity of $\mathbf{H}_{\mathbf{w}}$ and the triangle inequality, we have that:

$$\|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \mathbf{H}_{\mathbf{w}', \mathcal{K}, \mathcal{A}}\|_2 = \|\theta \mathbf{H}_{\mathbf{w}^*, \mathcal{K}} + (1 - \theta) \mathbf{H}_{\mathbf{w}^*, \mathcal{A}} - \theta \mathbf{H}_{\mathbf{w}', \mathcal{K}} - (1 - \theta) \mathbf{H}_{\mathbf{w}', \mathcal{A}}\|_2 \quad (\text{H.54})$$

$$\leq \theta \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}} - \mathbf{H}_{\mathbf{w}', \mathcal{K}}\|_2 + (1 - \theta) \|\mathbf{H}_{\mathbf{w}^*, \mathcal{A}} - \mathbf{H}_{\mathbf{w}', \mathcal{A}}\|_2 \quad (\text{H.55})$$

$$= \theta F_{\mathcal{K}} \|\mathbf{w}^* - \mathbf{w}'\|_2 + (1 - \theta) F_{\mathcal{A}} \|\mathbf{w}^* - \mathbf{w}'\|_2, \text{ by Lemma H.1} \quad (\text{H.56})$$

$$= \theta t F_{\mathcal{K}} \|\mathbf{w}^* - \tilde{\mathbf{w}}^*\|_2 + (1 - \theta) t F_{\mathcal{A}} \|\mathbf{w}^* - \tilde{\mathbf{w}}^*\|_2, \quad (\text{H.57})$$

This yields that:

$$\|\tilde{\mathbf{w}} - \tilde{\mathbf{w}}^*\|_2 \leq \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}^{-1}\|_2 \int_0^1 \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \mathbf{H}_{\mathbf{w}^* + t(\tilde{\mathbf{w}}^* - \mathbf{w}^*), \mathcal{K}, \mathcal{A}}\|_2 \|\mathbf{w}^* - \tilde{\mathbf{w}}^*\|_2 dt \quad (\text{H.58})$$

$$\leq \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}^{-1}\|_2 \int_0^1 (\theta t F_{\mathcal{K}} + (1 - \theta) t F_{\mathcal{A}}) \|\mathbf{w}^* - \tilde{\mathbf{w}}^*\|_2^2 dt \quad (\text{H.59})$$

$$= \frac{\theta F_{\mathcal{K}} + (1 - \theta) F_{\mathcal{A}}}{2} \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}^{-1}\|_2 \|\mathbf{w}^* - \tilde{\mathbf{w}}^*\|_2^2 \quad (\text{H.60})$$

as desired. \square

1674 H.6 PROOF OF PROPOSITION G.3
1675

1676 *Proof.* See the proof of theorem 3.4 in (Zhang et al., 2024), noting that in our setting $M = F =$
1677 $\theta F_{\mathcal{K}} + (1 - \theta)F_{\mathcal{A}}$ by Eq. (H.57). \square
1678

1679 H.7 PROOF OF THEOREM G.4
1680

1681 *Proof.* First, we have that:
1682

1683
1684
$$\mathbb{E}[\tilde{\mathbf{H}}_{t,\lambda}^{-1}] = \mathbb{E}[\mathbf{I} + \tilde{\mathbf{H}}_{t-1,\lambda}^{-1} - \frac{1}{J}\mathbf{H}_{t,\lambda}\tilde{\mathbf{H}}_{t-1,\lambda}^{-1}], \text{ by definition} \quad (\text{H.61})$$

1685
1686
$$= \mathbf{I} + \mathbb{E}[\tilde{\mathbf{H}}_{t-1,\lambda}^{-1}] - \frac{1}{J}\mathbb{E}[\mathbf{H}_{t,\lambda}\tilde{\mathbf{H}}_{t-1,\lambda}], \text{ linearity of expectation} \quad (\text{H.62})$$

1687
1688
$$= \mathbf{I} + \mathbb{E}[\tilde{\mathbf{H}}_{t-1,\lambda}^{-1}] - \frac{1}{J}\mathbb{E}[\mathbf{H}_{t,\lambda}]\mathbb{E}[\tilde{\mathbf{H}}_{t-1,\lambda}], \text{ i.i.d. samples} \quad (\text{H.63})$$

1689
1690
$$= \mathbf{I} + \mathbb{E}[\tilde{\mathbf{H}}_{t-1,\lambda}^{-1}] - \frac{\mathbf{H}_{w^*,\mathcal{K},\mathcal{A}} + \lambda\mathbf{I}}{J}\mathbb{E}[\tilde{\mathbf{H}}_{t-1,\lambda}], \text{ by Lemma H.5} \quad (\text{H.64})$$

1691

1692 Denote $\mathbf{H}_* := \mathbf{H}_{w^*,\mathcal{K},\mathcal{A}}$ and $\mathbf{E}_t := \mathbb{E}[\tilde{\mathbf{H}}_{t,\lambda}^{-1}]$. We thus have that:
1693

1694
1695
1696
$$\mathbf{E}_t = \mathbf{I} + \mathbf{E}_{t-1} - \frac{\mathbf{H}_*}{J}\mathbf{E}_{t-1} \quad (\text{H.65})$$

1697
1698
$$= \mathbf{I} + \mathbf{E}_{t-1}(\mathbf{I} - \frac{\mathbf{H}_*}{J}) \quad (\text{H.66})$$

1699
1700
$$= \mathbf{I} + (\mathbf{I} - \mathbf{M})\mathbf{E}_{t-1}, \text{ letting } \mathbf{M} := \frac{\mathbf{H}_*}{J} \quad (\text{H.67})$$

1701

1702 We then know that, by assumption, $\lambda > \|\mathbf{H}_*\|_2$, where \mathbf{H}_* is a symmetric Hessian by Lemma H.3;
1703 as such, $\mathbf{H}_* + \lambda\mathbf{I}$ is positive definite and has all positive eigenvalues. We also know that $\|\mathbf{H}_*\|_2 <$
1704 $J \implies \|\mathbf{M}\|_2 < 1$, so we have that $0 < \lambda_i(\mathbf{M}) < 1$ for all eigenvalues λ_i . Furthermore, $\mathbf{I} - \mathbf{M}$
1705 has eigenvalues $1 - \lambda_i(\mathbf{M})$, so we have that $0 < \lambda_i(\mathbf{I} - \mathbf{M}) < 1$, so $\|\mathbf{I} - \mathbf{M}\|_2 < 1$, since $\mathbf{I} - \mathbf{M}$
1706 is symmetric. Since $\mathbf{I} - \mathbf{M}$ has spectral radius less than 1, the Neumann series $\sum_{k=0}^{\infty}(\mathbf{I} - \mathbf{M})^k$
1707 converges. (Mayer, 1985). Thus, the Neumann series is Cauchy.

1708 Fix $\varepsilon > 0$. Let $s_n = \sum_{k=0}^n(\mathbf{I} - \mathbf{M})^k$. We know that $\exists N \in \mathbb{N}$ s.t. $m > n \geq N \implies$
1709 $\|s_m - s_n\|_2 = \|\sum_{k=n+1}^m(\mathbf{I} - \mathbf{M})^k\|_2 < \varepsilon$. For $m > n \geq N$, we have that $\|\mathbf{E}_m - \mathbf{E}_n\|_2 =$
1710 $\|\sum_{k=n+1}^m(\mathbf{I} - \mathbf{M})^k\|_2 < \varepsilon$. As such, $\{\mathbf{E}_n\}$ is Cauchy; since it is real, it converges. As such,
1711 $\mathbf{E}_{\infty} = \lim_{t \rightarrow \infty} \mathbf{E}_n$ exists.
1712

1713 Taking limits on both sides, we then have:
1714

1715
1716
$$\mathbb{E}[\tilde{\mathbf{H}}_{\infty,\lambda}^{-1}] = \mathbf{I} + \mathbb{E}[\tilde{\mathbf{H}}_{\infty,\lambda}^{-1}] + \frac{\mathbf{H}_{w^*,\mathcal{K},\mathcal{A}} + \lambda\mathbf{I}}{J}\mathbb{E}[\tilde{\mathbf{H}}_{\infty,\lambda}^{-1}] \quad (\text{H.68})$$

1717
1718
$$\iff \mathbb{E}[\frac{\tilde{\mathbf{H}}_{\infty,\lambda}^{-1}}{J}] = (\mathbf{H}_{w^*,\mathcal{K},\mathcal{A}} + \lambda\mathbf{I})^{-1} \quad (\text{H.69})$$

1719
1720

1721 rearranging using linearity of expectation and noting that λ was chosen such that $\mathbf{H}_{w^*,\mathcal{K},\mathcal{A}} + \lambda\mathbf{I}$ is
1722 invertible, as desired. \square
1723

1724 H.8 PROOF OF THEOREM G.5
1725

1726 This follows similarly to theorem 3.6 and proposition 4.1 in Zhang et al. (2024), noting that we
1727 apply Lemma H.7 instead of applying lemma 3.6 from Agarwal et al. (2016). Furthermore, note that
 $L = \theta P_{\mathcal{K}} + (1 - \theta)P_{\mathcal{A}}$ in our setting. For completeness, we provide the full proof below.

1728 *Proof.*

$$1729 \quad \tilde{\mathbf{w}} - \tilde{\mathbf{w}}^* = \mathbf{w}^* - \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H} \nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \tilde{\mathbf{w}}^* \quad (\text{H.70})$$

$$1731 \quad = \mathbf{w}^* - \tilde{\mathbf{w}}^* - \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H} (\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}}) - \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H} \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}} \quad (\text{H.71})$$

1732 By the triangle inequality, this yields:

$$1733 \quad \|\tilde{\mathbf{w}} - \tilde{\mathbf{w}}^*\|_2 \leq \|\mathbf{w}^* - \tilde{\mathbf{w}}^* - \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H} (\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}})\|_2 + \|\frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H} \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}}\|_2 \quad (\text{H.72})$$

1734 The first term in Eq. (H.72) can be bounded by the triangle inequality as:

$$1735 \quad \|\mathbf{w}^* - \tilde{\mathbf{w}}^* - \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H} (\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}})\|_2 \quad (\text{H.73})$$

$$1736 \quad = \|\mathbf{w}^* - \tilde{\mathbf{w}}^* - ((\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} + \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H} - (\mathbf{H}_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1}) (\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}})\|_2 \quad (\text{H.74})$$

$$1741 \quad \leq \|\tilde{\mathbf{w}}^* - \tilde{\mathbf{w}}^* - (\mathbf{H}_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} (\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}})\|_2 \quad (\text{H.75})$$

$$1742 \quad + \|((\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} - \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H}) (\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}})\|_2 \quad (\text{H.76})$$

1743 In the setting of Prop. G.3, we have that $\tilde{\mathbf{w}} - \tilde{\mathbf{w}}^* = \mathbf{w}^* - \tilde{\mathbf{w}}^* - (\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} (\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}})$. Hence, by Prop. G.3, we have that:

$$1744 \quad \|\mathbf{w}^* - \tilde{\mathbf{w}}^* - (\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} (\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}})\|_2 \quad (\text{H.77})$$

$$1745 \quad \leq \frac{2C((\theta F_{\mathcal{K}} + (1 - \theta) F_{\mathcal{A}})C + \lambda)}{\lambda + \lambda_{\min}} \quad (\text{H.78})$$

1746 Furthermore, we have:

$$1747 \quad \|((\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} - \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H}) (\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}})\|_2 \quad (\text{H.79})$$

$$1748 \quad \leq \|((\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} - \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H})\|_2 \|(\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}})\|_2, \text{property of op norm} \quad (\text{H.80})$$

$$1749 \quad \leq (16 \frac{B}{\zeta_{\min}} \sqrt{\frac{\ln \frac{d}{\rho}}{b} + \frac{1}{16}}) \|(\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}})\|_2, \text{Lemma H.7} \quad (\text{H.81})$$

$$1750 \quad \leq (16 \frac{B}{\zeta_{\min}} \sqrt{\frac{\ln \frac{d}{\rho}}{b} + \frac{1}{16}}) 2C(\theta P_{\mathcal{K}} + (1 - \theta) P_{\mathcal{A}}), \text{Lemma H.1} \quad (\text{H.82})$$

1751 with probability at least $1 - \rho$. Incorporating this into equation Eq. (H.76), we have that:

$$1752 \quad \|\mathbf{w}^* - \tilde{\mathbf{w}}^* - \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H} (\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}} - \nabla_{\tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}})\|_2 \leq \frac{2C((\theta F_{\mathcal{K}} + (1 - \theta) F_{\mathcal{A}})C + \lambda)}{\lambda + \lambda_{\min}} \quad (\text{H.83})$$

$$1753 \quad + (32 \frac{B}{\zeta_{\min}} \sqrt{\frac{\ln \frac{d}{\rho}}{b} + \frac{1}{8}}) C(\theta P_{\mathcal{K}} + (1 - \theta) P_{\mathcal{A}}) \quad (\text{H.84})$$

1782 It then suffices to bound the second term in Eq. (H.72). We have that:

$$1783 \left\| \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H} \nabla \tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A} \right\|_2 = \left\| \left[(\mathbf{H}_{w^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} - (\tilde{\mathbf{H}}_{n,\lambda}^{-1})^{-1} \right] \nabla \tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A} \right\|_2 \quad (\text{H.85})$$

$$1787 = \left\| (\mathbf{H}_{w^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} \nabla \tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A} \right. \\ 1788 \left. + \left(\frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H} - (\mathbf{H}_{w^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} \right) \nabla \tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A} \right\|_2 \quad (\text{H.86})$$

$$1791 \leq \|(\mathbf{H}_{w^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1}\|_2 \|\nabla \tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}\|_2 \\ 1792 + \left\| \frac{\tilde{\mathbf{H}}_{n,\lambda}^{-1}}{H} - (\mathbf{H}_{w^*, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1} \right\|_2 \|\nabla \tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}\|_2 \quad (\text{H.87})$$

$$1793 \leq \frac{G}{\lambda + \lambda_{\min}} + \left(16 \frac{B}{\zeta_{\min}} \sqrt{\frac{\ln(d/\rho)}{b}} + \frac{1}{16} \right) G. \quad (\text{H.88})$$

1796 by definition of λ , Lemma H.7, and that $\|\nabla \tilde{\mathbf{w}}^*, \mathcal{K}, \mathcal{A}\|_2 \leq G$.

1798 Incorporating the above into Eq. (H.72), this yields that:

$$1800 \|\tilde{\mathbf{w}} - \tilde{\mathbf{w}}^*\|_2 \leq \frac{2C((\theta F_{\mathcal{K}} + (1 - \theta)F_{\mathcal{A}})C + \lambda)}{\lambda + \lambda_{\min}} \quad (\text{H.89})$$

$$1803 + \left(32 \frac{B}{\zeta_{\min}} \sqrt{\frac{\ln(d/\rho)}{b}} + \frac{1}{8} \right) C \left(\theta P_{\mathcal{K}} + (1 - \theta)P_{\mathcal{A}} + \frac{G}{\lambda + \lambda_{\min}} \right. \\ 1804 \left. + \left(16 \frac{B}{\zeta_{\min}} \sqrt{\frac{\ln(d/\rho)}{b}} + \frac{1}{16} \right) G \right) \quad (\text{H.90})$$

$$1808 = \frac{2C((\theta F_{\mathcal{K}} + (1 - \theta)F_{\mathcal{A}})C + \mu) + G}{\mu + \mu_{\min}} \quad (\text{H.91}) \\ 1810 + \left(16 \frac{B}{\zeta_{\min}} \sqrt{\frac{\ln(d/\rho)}{b}} + \frac{1}{16} \right) (2C(\theta P_{\mathcal{K}} + (1 - \theta)P_{\mathcal{A}} + G)).$$

1813 as desired.

1814 \square

1816 H.9 PROOF OF PROPOSITION I.1

1818 *Proof.* By the same token as proposition 4.2 in Zhang et al. (2024), follow the proof of theorem G.5.

1819 \square

1821 H.10 PROOF OF PROPOSITION 5.1

1823 *Proof.* Fix any sampled \mathcal{D} . Since $\mathcal{M}_\theta(\mathcal{D})$ is taken to be the global risk minimizer, we have that:

$$1825 \theta \mathcal{L}_{\mathcal{K}}(\mathcal{M}_\theta(\mathcal{D}), \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathcal{M}_\theta(\mathcal{D}), \mathcal{D}_r) \quad (\text{H.92})$$

$$1826 \leq \theta \mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r) \quad \forall \mathbf{w} \in \mathcal{W} \quad (\text{H.93})$$

1828 subtracting $\frac{\lambda}{2} \|\mathbf{w}\|_2$ from both sides.

1830 Let \mathbf{w}_U be the parameter that results in a parameterized model $f_{\mathbf{w}_U}$ which outputs a uniform distribution; by Prop. 4.2, such a parameter exists. We then have that:

$$1833 \mathcal{L}_{\mathcal{K}}(\mathbf{w}_U, \mathcal{D}_f) = \sum_{i=1}^{|\mathcal{D}_f|} D_{KL}(U[0, |\mathcal{Y}|]) \|U[0, |\mathcal{Y}|] = 0 \quad (\text{H.94})$$

1836 and

$$1837 \mathcal{L}_{\mathcal{A}}(\mathbf{w}_U, \mathcal{D}_r) = \sum_{i=1}^{|\mathcal{D}_r|} \mathbb{H}_{CE}(\mathbf{y}^{(i)}, U[0, |\mathcal{Y}|]) = - \sum_{i=1}^{|\mathcal{D}_r|} \sum_{j=1}^{|\mathcal{Y}|} \mathbf{y}_j^{(i)} \ln \frac{1}{|\mathcal{Y}|} = |\mathcal{D}_r| \ln |\mathcal{Y}| \quad (\text{H.95})$$

1842 where \mathbf{y} is a one hot vector of length \mathcal{Y} such that for $\mathbf{y}_j^{(i)}$, $j = 1, \dots, |\mathcal{Y}|$,

$$1843 \mathbf{y}_j^{(i)} = \begin{cases} 1 & \text{instance } i \text{ is labeled class } j \\ 0 & \text{instance } i \text{ is not labeled class } j \end{cases} \quad (\text{H.96})$$

1844 Incorporating the above into Eq. (H.93) yields:

$$1845 \theta \mathcal{L}_{\mathcal{K}}(\mathcal{M}_{\theta}(\mathcal{D}), \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathcal{M}_{\theta}(\mathcal{D}), \mathcal{D}_r) \leq \theta \mathcal{L}_{\mathcal{K}}(\mathbf{w}_U, \mathcal{D}_f) + (1 - \theta) \mathcal{L}_{\mathcal{A}}(\mathbf{w}_U, \mathcal{D}_r) \quad (\text{H.97})$$

$$1846 \leq \theta(0) + (1 - \theta) |\mathcal{D}_r| \ln |\mathcal{Y}| \quad (\text{H.98})$$

$$1847 = |\mathcal{D}_r| (1 - \theta) \ln |\mathcal{Y}| \quad (\text{H.99})$$

1848 This then yields that:

$$1849 \mathcal{L}_{\mathcal{K}}(\mathcal{M}_{\theta}(\mathcal{D}), \mathcal{D}_f) \leq \frac{1 - \theta}{\theta} (|\mathcal{D}_r| \ln |\mathcal{Y}| - \mathcal{L}_{\mathcal{A}}(\mathcal{M}_{\theta}(\mathcal{D}), \mathcal{D}_r)) \quad (\text{H.100})$$

$$1850 \leq \frac{1 - \theta}{\theta} |\mathcal{D}_r| \ln |\mathcal{Y}| \quad (\text{H.101})$$

1851 since the cross entropy is nonnegative, yielding that $-\mathcal{L}_{\mathcal{A}}(\mathcal{M}_{\theta}(\mathcal{D}), \mathcal{D}_r) \leq 0$.

1852 Then, we have:

$$1853 \|f_{\mathcal{M}_{\theta}(\mathcal{D})}(\mathcal{D}_f) - U[0, |\mathcal{Y}|]\|_{\infty} \leq \|f_{\mathcal{M}_{\theta}(\mathcal{D})}(\mathcal{D}_f) - U[0, |\mathcal{Y}|]\|_1 \quad (\text{H.102})$$

$$1854 \leq 2TV(f_{\mathcal{M}_{\theta}(\mathcal{D})}(\mathcal{D}_f) - U[0, |\mathcal{Y}|]) \quad (\text{H.103})$$

$$1855 \leq 2\sqrt{\frac{1}{2} D_{KL}(f_{\mathcal{M}_{\theta}(\mathcal{D})}, \|U[0, |\mathcal{Y}]\|)}, \text{ Pinsker's inequality (Pinsker, 1964)}$$

$$1856 \quad (\text{H.104})$$

$$1857 = \sqrt{2D_{KL}(f_{\mathcal{M}_{\theta}(\mathcal{D})}, \|U[0, |\mathcal{Y}]\|)} \quad (\text{H.105})$$

$$1858 = \sqrt{2\mathcal{L}_{\mathcal{K}}(\mathcal{M}_{\theta}(\mathcal{D}), \mathcal{D}_f)} \quad (\text{H.106})$$

$$1859 \leq \sqrt{2|\mathcal{D}_r| \left(\frac{1 - \theta}{\theta}\right) \ln |\mathcal{Y}|} \quad (\text{H.107})$$

1860 by the above bound on $\mathcal{L}_{\mathcal{K}}$, as desired. \square

1861 H.11 PROOF OF COROLLARY 5.2

1862 *Proof.* To have Eq. (7), by Prop. 5.1, it suffices to solve for θ in the bound obtained. This results in:

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902

$$\sqrt{2\left(\frac{1-\theta}{\theta}\right)|\mathcal{D}_r|\ln|\mathcal{Y}|} \leq \varepsilon \iff \frac{1-\theta}{\theta}|\mathcal{D}_r|\ln|\mathcal{Y}| \leq \frac{\varepsilon^2}{2} \quad (\text{H.108})$$

$$\iff \frac{|\mathcal{D}_r|\ln|\mathcal{Y}|}{\theta} - \frac{|\mathcal{D}_r|\ln|\mathcal{Y}|\theta}{\theta} \leq \frac{\varepsilon^2}{2} \quad (\text{H.109})$$

$$\iff \frac{|\mathcal{D}_r|\ln|\mathcal{Y}|}{\theta} \leq \frac{\varepsilon^2}{2} + |\mathcal{D}_r|\ln|\mathcal{Y}| = \frac{\varepsilon^2 + 2|\mathcal{D}_r|\ln|\mathcal{Y}|}{2} \quad (\text{H.110})$$

$$\iff \frac{\theta}{|\mathcal{D}_r|\ln|\mathcal{Y}|} \geq \frac{2}{\varepsilon^2 + 2|\mathcal{D}_r|\ln|\mathcal{Y}|} \quad (\text{H.111})$$

$$\iff \theta \geq \frac{2|\mathcal{D}_r|\ln|\mathcal{Y}|}{\varepsilon^2 + 2|\mathcal{D}_r|\ln|\mathcal{Y}|} \quad (\text{H.112})$$

1903 as desired.

1904
1905 \square

1906 H.12 PROOF OF THEOREM 5.5

1907 First, before we prove theorem 5.5, we note that we can use Lemma H.2 and that $\|\mathbf{w}\| \leq 2$ to obtain
1908 a simple bound. Let:

$$1909 |\alpha^* - \alpha(\theta)| = |\mathcal{L}_{\mathcal{A}}(\mathcal{A}(\mathcal{D}_r), \mathcal{D}_r) - \mathcal{L}_{\mathcal{A}}(\mathcal{M}_{\theta}(\mathcal{D}), \mathcal{D}_r)| \quad (\text{H.113})$$

$$1910 \leq \frac{P_{\mathcal{A}}}{2} \|\mathcal{M}_{\theta}(\mathcal{D}) - \mathcal{A}(\mathcal{D}_r)\|_2^2 + \|\nabla_{\mathcal{A}(\mathcal{D}_r), \mathcal{A}}\|_2 \|\mathcal{M}_{\theta}(\mathcal{D}) - \mathcal{A}(\mathcal{D}_r)\|_2 \quad (\text{H.114})$$

$$1911 \leq \frac{C^2 P_{\mathcal{A}}}{2} + \lambda C^2 \quad (\text{H.115})$$

1912 after applying the triangle inequality and rearranging the first order condition on $\mathcal{A}(\mathcal{D}_r)$.

1913 However, this bound is vacuous and not tight; it does not incorporate any information about θ or most
1914 of the constants that appear in Asm. 5.3 and Asm. 5.4. Given this, we seek to construct a tighter,
1915 non-vacuous bound. We first restate the proof without any asymptotic characterizations:

1916 **Theorem H.11.** *Suppose Assumptions 5.3 and 5.4 hold, and let $P_{\mathcal{K}}, P_{\mathcal{K}}, F_{\mathcal{K}}, F_{\mathcal{A}}$ be as defined in
1917 Assumptions 5.3 and 5.4. Let $\alpha^* := \mathcal{L}_{\mathcal{A}}(\mathcal{A}(\mathcal{D}_r), \mathcal{D}_r)$ be the locally optimal (empirical) retain loss,
1918 achieved by $\mathcal{M}_{\theta}(\mathcal{D})$ when $\theta = 0$. Let $\alpha(\theta) := \mathcal{L}_{\mathcal{A}}(\mathcal{M}_{\theta}(\mathcal{D}), \mathcal{D}_r)$ be the locally optimal retain
1919 loss obtained by $\mathcal{M}_{\theta}(\mathcal{D})$ when $\theta \in (0, 1)$. Suppose all weights used throughout are bounded by
1920 $\|\mathbf{w}\|_2 \leq C$. Additionally, denote by $F := \theta M_{\mathcal{K}} + (1 - \theta)F_{\mathcal{A}}$ and $P := \theta P_{\mathcal{K}} + (1 - \theta)P_{\mathcal{A}}$.
1921 Consider regularization coefficient $\lambda \geq L + 2\theta CF + \sqrt{2\theta CF(P + 2\theta CF + 8P_{\mathcal{K}})}$. Then, we have
1922 the following bound:*

$$1923 |\alpha^* - \alpha(\theta)| \leq \frac{P_{\mathcal{K}}}{2} \left(\frac{\lambda - P - \sqrt{(\lambda - P)^2 - 4\theta CF(2P_{\mathcal{K}} + \lambda)}}{2F} \right)^2 + \quad (\text{H.116})$$

$$1924 \lambda C \left(\frac{\lambda - P - \sqrt{(\lambda - P)^2 - 4\theta CF(2P_{\mathcal{K}} + \lambda)}}{2F} \right). \quad (\text{H.117})$$

1925 *Proof.* First, when $\theta = 0$, we have that:

$$1926 \mathbf{w}_{\alpha^*} := \arg \min_{\mathbf{w} \in \mathcal{W}, \|\mathbf{w}\|_2 \leq C} \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (\text{H.118})$$

1927 which yields the first order condition:

$$1928 \nabla_{\mathbf{w}_{\alpha^*}, \mathcal{A}} + \lambda \mathbf{w}_{\alpha^*} = 0 \quad (\text{H.119})$$

1944 which, upon multiplying $1 - \theta$ on both sides, yields:
1945

$$1946 \quad (1 - \theta)\nabla_{\mathbf{w}_{\alpha^*}, \mathcal{A}} + (1 - \theta)\lambda\mathbf{w}_{\alpha^*} = 0 \quad (\text{H.120})$$

1948 Then, when $\theta \in (0, 1)$, we have:
1949

$$1950 \quad \mathbf{w}_{\alpha(\theta)} := \arg \min_{\mathbf{w} \in \mathcal{W}, \|\mathbf{w}\|_2 \leq C} \theta \mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f) + (1 - \theta)\mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r) + \frac{\lambda}{2}\|\mathbf{w}\|_2^2 \quad (\text{H.121})$$

1953 which yields the first order condition:
1954

$$1955 \quad \theta \nabla_{\mathbf{w}_{\alpha(\theta)}, \mathcal{K}} + (1 - \theta)\nabla_{\mathbf{w}_{\alpha(\theta)}, \mathcal{A}} + \lambda\mathbf{w}_{\alpha(\theta)} = 0 \quad (\text{H.122})$$

1957 Subtracting Eq. (H.120) from Eq. (H.122) yields:
1958

$$1959 \quad \theta \nabla_{\mathbf{w}_{\alpha(\theta)}, \mathcal{K}} + (1 - \theta)\nabla_{\mathbf{w}_{\alpha(\theta)}, \mathcal{A}} + \lambda\mathbf{w}_{\alpha(\theta)} - (1 - \theta)\nabla_{\mathbf{w}_{\alpha^*}, \mathcal{A}} - (1 - \theta)\lambda\mathbf{w}_{\alpha^*} = 0 \quad (\text{H.123})$$

1961 which simplifies to:
1962

$$1963 \quad \theta \nabla_{\mathbf{w}_{\alpha(\theta)}, \mathcal{K}} + (1 - \theta)(\nabla_{\mathbf{w}_{\alpha(\theta)}, \mathcal{A}} - \nabla_{\mathbf{w}_{\alpha^*}, \mathcal{A}}) + \lambda(\mathbf{w}_{\alpha(\theta)} - \mathbf{w}_{\alpha^*}) = -\theta\lambda\mathbf{w}_{\alpha^*} \quad (\text{H.124})$$

1965 The fundamental theorem of calculus then yields:
1966

$$1967 \quad \int_0^1 \mathbf{H}_{\mathbf{w}_{\alpha^*} + t(\mathbf{w}_{\alpha(\theta)} + \mathbf{w}_{\alpha^*}), \mathcal{A}}(\mathbf{w}_{\alpha(\theta)} - \mathbf{w}_{\alpha^*}) dt = \nabla_{\mathbf{w}_{\alpha(\theta)}, \mathcal{A}} - \nabla_{\mathbf{w}_{\alpha^*}, \mathcal{A}} \quad (\text{H.125})$$

1970 and
1971

$$1972 \quad \int_0^1 \mathbf{H}_{\mathbf{w}_{\alpha^*} + t(\mathbf{w}_{\alpha(\theta)} + \mathbf{w}_{\alpha^*}), \mathcal{K}}(\mathbf{w}_{\alpha(\theta)} - \mathbf{w}_{\alpha^*}) dt = \nabla_{\mathbf{w}_{\alpha(\theta)}, \mathcal{K}} - \nabla_{\mathbf{w}_{\alpha^*}, \mathcal{K}} \quad (\text{H.126})$$

1975 We thus denote:
1976

$$1977 \quad \bar{\mathbf{H}}_{\mathcal{K}} := \int_0^1 \mathbf{H}_{\mathbf{w}_{\alpha^*} + t(\mathbf{w}_{\alpha(\theta)} + \mathbf{w}_{\alpha^*}), \mathcal{K}} dt \quad (\text{H.127})$$

$$1980 \quad \bar{\mathbf{H}}_{\mathcal{A}} := \int_0^1 \mathbf{H}_{\mathbf{w}_{\alpha^*} + t(\mathbf{w}_{\alpha(\theta)} + \mathbf{w}_{\alpha^*}), \mathcal{A}} dt \quad (\text{H.128})$$

$$1983 \quad \Delta\mathbf{w} := \mathbf{w}_{\alpha(\theta)} - \mathbf{w}_{\alpha^*} \quad (\text{H.129})$$

1984
1985 Incorporating Eq. (H.125) and Eq. (H.126) into Eq. (H.124) then yields:
1986

$$1987 \quad \begin{aligned} & \theta(\nabla_{\mathbf{w}_{\alpha^*}, \mathcal{K}} + \bar{\mathbf{H}}_{\mathcal{K}}\Delta\mathbf{w}) + (1 - \theta)(\nabla_{\mathbf{w}_{\alpha^*}, \mathcal{A}} + \bar{\mathbf{H}}_{\mathcal{A}}\Delta\mathbf{w} - \nabla_{\mathbf{w}_{\alpha^*}, \mathcal{A}}) + \lambda\Delta\mathbf{w} = -\theta\lambda\mathbf{w}_{\alpha^*} \\ \iff & \theta \nabla_{\mathbf{w}_{\alpha^*}, \mathcal{K}} + \theta \bar{\mathbf{H}}_{\mathcal{K}}\Delta\mathbf{w} + (1 - \theta)\bar{\mathbf{H}}_{\mathcal{A}}\Delta\mathbf{w} + \lambda\Delta\mathbf{w} + \theta\lambda\mathbf{w}_{\alpha^*} = 0 \\ \iff & (\theta \bar{\mathbf{H}}_{\mathcal{K}} + (1 - \theta)\bar{\mathbf{H}}_{\mathcal{A}} + \lambda\mathbf{I})\Delta\mathbf{w} = -\theta(\nabla_{\mathbf{w}_{\alpha^*}, \mathcal{K}} + \lambda\mathbf{w}_{\alpha^*}) \\ \iff & (\mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}, \mathcal{A}} + \lambda\mathbf{I} + \theta(\bar{\mathbf{H}}_{\mathcal{K}} - \mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}}) + (1 - \theta)(\bar{\mathbf{H}}_{\mathcal{A}} - \mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{A}}))\Delta\mathbf{w} \\ & = -\theta(\nabla_{\mathbf{w}_{\alpha^*}, \mathcal{K}} + \lambda\mathbf{w}_{\alpha^*}) \\ \iff & (\mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}, \mathcal{A}} + \lambda\mathbf{I})\Delta\mathbf{w} = -(\theta(\bar{\mathbf{H}}_{\mathcal{K}} - \mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}}) + (1 - \theta)(\bar{\mathbf{H}}_{\mathcal{A}} - \mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{A}}))\Delta\mathbf{w} \\ & \quad - \theta(\nabla_{\mathbf{w}_{\alpha^*}, \mathcal{K}} + \lambda\mathbf{w}_{\alpha^*}). \end{aligned} \quad (\text{H.130})$$

1997 Then, note that:

1998
 1999
 2000
 2001
 2002
 2003
 2004
 2005
 2006
 2007
 2008
 2009
 2010
 2011
 2012
 2013
 2014
 2015
 2016
 2017
 2018
 2019
 2020
 2021
 2022
 2023
 2024
 2025
 2026
 2027
 2028
 2029
 2030
 2031
 2032
 2033
 2034
 2035
 2036
 2037
 2038
 2039
 2040
 2041
 2042
 2043
 2044
 2045
 2046
 2047
 2048
 2049
 2050
 2051

$$\|\bar{\mathbf{H}}_{\mathcal{K}} - \mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}}\|_2 = \left\| \int_0^1 \mathbf{H}_{\mathbf{w}_{\alpha^*} + t\Delta\mathbf{w}, \mathcal{K}} dt - \mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}} \right\|_2 \quad (\text{H.131})$$

$$\leq \int_0^1 \|\mathbf{H}_{\mathbf{w}_{\alpha^*} + t\Delta\mathbf{w}, \mathcal{K}} - \mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}}\|_2 dt \quad (\text{H.132})$$

$$\leq \frac{F_{\mathcal{K}}}{2} \|\Delta\mathbf{w}\|_2 \quad (\text{H.133})$$

by the same token as in Prop. G.3.

Similarly:

$$\|\bar{\mathbf{H}}_{\mathcal{A}} - \mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{A}}\|_2 \leq \frac{F_{\mathcal{A}}}{2} \|\Delta\mathbf{w}\|_2 \quad (\text{H.134})$$

Also:

$$\|\nabla_{\mathbf{w}_{\alpha^*}, \mathcal{K}}\|_2 = \|\nabla_{\mathbf{w}_{\alpha^*}, \mathcal{K}} - \nabla_{\mathcal{K}(\mathcal{D}_f), \mathcal{K}}\|_2 \quad (\text{H.135})$$

$$\leq P_{\mathcal{K}} \|\mathbf{w}_{\alpha^*} - \mathbf{w}_{\mathcal{K}(\mathcal{D}_f)}\|_2 \quad (\text{H.136})$$

$$\leq 2P_{\mathcal{K}}C \quad (\text{H.137})$$

by definition of the uniform learner \mathcal{K} , Lemma H.1, and the triangle inequality.

Additionally:

$$\|\lambda\mathbf{w}_{\alpha^*}\|_2 \leq \lambda C \quad (\text{H.138})$$

By the triangle inequality, incorporating Eq. (H.133), Eq. (H.134), Eq. (H.137), and Eq. (H.138) into Eq. (H.130), we have that:

$$\|(\mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}, \mathcal{A}} + \lambda\mathbf{I})\Delta\mathbf{w}\|_2 \leq (\theta\|\bar{\mathbf{H}}_{\mathcal{K}} - \mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}}\|_2 + (1-\theta)\|\bar{\mathbf{H}}_{\mathcal{A}} - \mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{A}}\|_2) \|\Delta\mathbf{w}\|_2 + \theta\|\nabla_{\mathbf{w}_{\alpha^*}, \mathcal{K}}\|_2 + \theta\|\lambda\mathbf{w}_{\alpha^*}\|_2 \quad (\text{H.139})$$

$$\leq (\theta F_{\mathcal{K}} + (1-\theta)F_{\mathcal{A}}) \|\Delta\mathbf{w}\|_2 + \theta C(2P_{\mathcal{K}} + \lambda) \|\Delta\mathbf{w}\|_2 \quad (\text{H.140})$$

$$\leq (\theta F_{\mathcal{K}} + (1-\theta)F_{\mathcal{A}}) \|\Delta\mathbf{w}\|_2 + \theta C(2P_{\mathcal{K}} + \lambda) \|\Delta\mathbf{w}\|_2 \quad (\text{H.141})$$

Note that we have, where $\sigma_{\min}(\cdot)$ denotes the minimum singular value:

$$\|(\mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}, \mathcal{A}} + \lambda\mathbf{I})\Delta\mathbf{w}\|_2 \geq \sigma_{\min}(\mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}, \mathcal{A}} + \lambda\mathbf{I}) \|\Delta\mathbf{w}\|_2 \text{ by property of op. norm} \quad (\text{H.142})$$

$$= \lambda_{\min}(\mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}, \mathcal{A}} + \lambda\mathbf{I}) \|\Delta\mathbf{w}\|_2 \text{ by Lemma H.3} \quad (\text{H.143})$$

Furthermore, by Lemma H.1, we have that $\|\mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}}\|_2 \leq P_{\mathcal{K}}$ and $\|\mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{A}}\|_2 \leq P_{\mathcal{A}}$, which yields:

$$\|\mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}, \mathcal{A}}\|_2 \leq \theta P_{\mathcal{K}} + (1-\theta)P_{\mathcal{A}} \quad (\text{H.144})$$

which by Lemma H.3 yields:

$$\lambda_{\min}(\mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}, \mathcal{A}} + \lambda\mathbf{I}) \in [\lambda - \theta P_{\mathcal{K}} - (1-\theta)P_{\mathcal{A}}, \mu + \theta P_{\mathcal{K}} + (1-\theta)P_{\mathcal{A}}] \quad (\text{H.145})$$

With Eq. (H.143), this yields that:

$$\|(\mathbf{H}_{\mathbf{w}_{\alpha^*}, \mathcal{K}, \mathcal{A}} + \lambda\mathbf{I})\Delta\mathbf{w}\|_2 \geq (\lambda - \theta P_{\mathcal{K}} - (1-\theta)P_{\mathcal{A}}) \|\Delta\mathbf{w}\|_2 \quad (\text{H.146})$$

2052 Incorporating this into Eq. (H.141) yields:
 2053
 2054

$$2055 \quad (\lambda - \theta P_{\mathcal{K}} - (1 - \theta)P_{\mathcal{A}})\|\Delta\mathbf{w}\|_2 \leq (\theta F_{\mathcal{K}} + (1 - \theta)F_{\mathcal{A}})\|\Delta\mathbf{w}\|_2^2 + \theta C(2P_{\mathcal{K}} + \lambda) \quad (\text{H.147})$$

2057 Simplifying yields the quadratic inequality:
 2058
 2059

$$2060 \quad (\theta F_{\mathcal{K}} + (1 - \theta)F_{\mathcal{A}})\|\Delta\mathbf{w}\|_2^2 - (\lambda - \theta P_{\mathcal{K}} - (1 - \theta)P_{\mathcal{A}})\|\Delta\mathbf{w}\|_2 + \theta C(2P_{\mathcal{K}} + \lambda) \geq 0 \quad (\text{H.148})$$

2062 This then yields that:
 2063
 2064

$$2065 \quad \|\Delta\mathbf{w}\|_2 \leq \frac{\lambda - P - \sqrt{(\lambda - P)^2 - 4\theta CF(2P_{\mathcal{K}} + \lambda)}}{2F} \quad (\text{H.149})$$

2068 This is only valid when:
 2069

$$2070 \quad (\lambda - P)^2 - 4\theta CF(2P_{\mathcal{K}} + \lambda) \geq 0 \quad (\text{H.150})$$

$$2071 \quad \iff \lambda^2 - 2L\lambda + P^2 - 4\theta C2P_{\mathcal{K}}F - 4\theta C\lambda F \geq 0 \quad (\text{H.151})$$

$$2072 \quad \iff \lambda^2 - 2P\lambda - 4\theta CF\lambda + P^2 - 8\theta CP_{\mathcal{K}}F \geq 0 \quad (\text{H.152})$$

$$2073 \quad \iff \lambda^2 - (2P + 4\theta CF)\lambda + (P^2 - 8\theta CP_{\mathcal{K}}F) \geq 0 \quad (\text{H.153})$$

$$2074 \quad \iff \lambda \geq P + 2\theta CF + \sqrt{2\theta CF(P + 2\theta CF + 8P_{\mathcal{K}})} \quad (\text{H.154})$$

2077 which holds by assumption. Note that all components of $2\theta CF + \sqrt{2\theta CF(P + 2\theta CF + 8P_{\mathcal{K}})}$ are
 2078 nonnegative, rendering this valid. Incorporating Eq. (H.149) into Lemma H.2 yields the final bound
 2079 as desired. \square
 2080

2081 Then, theorem 5.5 follows as a corollary of theorem H.11:
 2082

2083 *Proof.* Note that we take care to ensure the bound holds for any choice of $\theta \in [0, 1]$. Hence, fix
 2084 $\theta \in [0, 1]$.
 2085

2086 Let

$$2087 \quad a := \lambda - P > 0, \quad \varepsilon := 4\theta CF(2P_{\mathcal{K}} + \lambda), \quad \Delta := \frac{a - \sqrt{a^2 - \varepsilon}}{2F}. \quad (\text{H.155})$$

2089 Theorem H.11 gives the inequality:
 2090

$$2091 \quad |\alpha^* - \alpha(\theta)| \leq \frac{P_{\mathcal{K}}}{2} \Delta^2 + \lambda C \Delta. \quad (\text{H.156})$$

2094 By the condition on λ in the theorem, the square root is real for every $\theta \in [0, 1]$ (i.e. $a^2 - \varepsilon \geq 0$).
 2095 This yields:
 2096

$$2097 \quad a - \sqrt{a^2 - \varepsilon} = \frac{\varepsilon}{a + \sqrt{a^2 - \varepsilon}}. \quad (\text{H.157})$$

2100 Then, since $a + \sqrt{a^2 - \varepsilon} \geq a > 0$, (H.157) implies:
 2101

$$2102 \quad a - \sqrt{a^2 - \varepsilon} \leq \frac{\varepsilon}{a}. \quad (\text{H.158})$$

2104 Dividing (H.158) by $2F$ yields:
 2105

$$\Delta \leq \frac{\varepsilon}{2aF}. \quad (\text{H.159})$$

Then, substituting $\varepsilon = 4\theta CF(2P_K + \lambda)$ from (H.155) yields:

$$\Delta \leq \frac{4\theta CF(2P_K + \lambda)}{2aF} = \frac{2\theta C(2P_K + \lambda)}{a}. \quad (\text{H.160})$$

We now bound the two terms on the right-hand side of (H.156). Using (H.160), we have that:

$$\lambda C \Delta \leq \lambda C \cdot \frac{2\theta C(2P_K + \lambda)}{a} = \frac{2\lambda(2P_K + \lambda)}{a} C^2 \theta = \mathcal{O}(\lambda C^2 \theta), \quad (\text{H.161})$$

$$\frac{P_K}{2} \Delta^2 \leq \frac{P_K}{2} \left(\frac{2\theta C(2P_K + \lambda)}{a} \right)^2 = \frac{2P_K(2P_K + \lambda)^2}{a^2} C^2 \theta^2 = \mathcal{O}(C^2 \theta^2). \quad (\text{H.162})$$

Combining (H.156), (H.161) and (H.162) and absorbing constants (which are independent of $\theta \in [0, 1]$) yields:

$$|\alpha^* - \alpha(\theta)| = \mathcal{O}(\lambda C^2 \theta + C^2 \theta^2), \quad \text{for any } \theta \in [0, 1]. \quad (\text{H.163})$$

as desired. \square

I ONLINE ALGORITHM

We also consider the online setting, where users send requests in sequential order (Nguyen et al., 2022). Here, we denote \mathcal{D}_{f_k} as the forget set after the k -th request and the associated retain set as $\mathcal{D}_{r_k} = \mathcal{D} \setminus \cup_{i=1}^k \mathcal{D}_{f_i}$. Letting $\tilde{\mathbf{w}}_0 = \mathcal{A}(\mathcal{D})$, we estimate $\tilde{\mathbf{w}}_k$ recursively as $\tilde{\mathbf{w}}_k = \tilde{\mathbf{w}}_{k-1} - \frac{\tilde{H}_{n,\lambda,k-1}^{-1}}{H} \nabla_{\tilde{\mathbf{w}}_{k-1}, \mathcal{K}, \mathcal{A}}$, where $\tilde{H}_{n,\lambda,k-1}^{-1}$ is an estimator for $(\mathbf{H}_{\tilde{\mathbf{w}}_{k-1}, \mathcal{K}, \mathcal{A}} + \lambda \mathbf{I})^{-1}$ with respect to \mathcal{D}_{f_k} and \mathcal{D}_{r_k} . $\nabla_{\tilde{\mathbf{w}}_{k-1}, \mathcal{K}, \mathcal{A}}$ is also computed with respect to \mathcal{D}_{f_k} and \mathcal{D}_{r_k} . Adding noise to $\tilde{\mathbf{w}}_k$ as stipulated in theorem G.1 yields a \mathbf{w}_k^- satisfying Def. 4.6. Furthermore, we have that:

Proposition I.1. *Let λ_{\min} be the smallest eigenvalue of $\mathbf{H}_{\tilde{\mathbf{w}}_{k-1}, \mathcal{K}, \mathcal{A}}$, $\lambda > \|\mathbf{H}_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}\|_2$, and $\|\nabla_{\tilde{\mathbf{w}}_k, \mathcal{K}, \mathcal{A}}\|_2, \|\nabla_{\tilde{\mathbf{w}}_k, \mathcal{K}, \mathcal{A}}\|_2 \leq G$ for all k , all evaluated with respect to \mathcal{D}_{f_k} and \mathcal{D}_{r_k} . Then, the bound in theorem G.5 is identical in the online setting.*

Proof: See Appx. H.9.

In the online setting, Prop. I.1 yields Alg. 4.

Note that, for simplicity, we set $b = 1$. However, they can be added similarly to Alg. 3 if the user desires.

J ELIMINATING HYPERPARAMETERS IN CERTIFIED ALGORITHMS

Here, we summarize how to eliminate hyperparameters in Alg. 2, Alg. 3, and Alg. 4.

- λ_{\min} can be chosen as 0 by convex approximation, or it can be estimated using simple algorithms like Gershgorin's circle theorem or inverse power iteration.
- By Lemma H.1, λ can be chosen as $\theta P_{\mathcal{K}} + (1 - \theta) P_{\mathcal{A}}$
- Similarly, by Lemma H.1, H can be chosen as 2λ
- In practice, since $\frac{B}{\lambda + \lambda_{\min}}$ offers a bound on $\hat{\kappa}_l$ by Lemma H.6, ζ_{\min} can be chosen as $\lambda + \lambda_{\min}$
- By Lemma H.7, n can be chosen as $n = 2 \frac{B}{\lambda + \lambda_{\min}} \ln(\frac{B}{\lambda + \lambda_{\min}} b)$
- G can be approximated by computing $\|\nabla_{\mathbf{w}^*, \mathcal{K}, \mathcal{A}}\|_2$.
- θ can be chosen with Cor. 5.2 to satisfy a particular closeness to uniformity.
- In practice, we find that C can be chosen as 10, 20, or 100.
- b can be chosen to satisfy a particular concentration on the estimator, so we let it be free. However, one can set $b = 1$.
- Common heuristics for ε and δ are available in the differential privacy literature.
- In practice, following what is common in certified unlearning e.g. in (Zhang et al., 2024), the Lipschitz constants in Asm. 5.3 and Asm. 5.4 are treated as hyperparameters. However, in practice, they can all be set to 1.

Algorithm 4 Online $(\varepsilon, \delta, \theta)$ -certified uniformity with DP

Require: Dataset \mathcal{D} ; forget sets $\{\mathcal{D}_{f_1}, \dots, \mathcal{D}_{f_k}\}$; pretrained model $\mathbf{w}^* = \mathcal{A}(\mathcal{D})$; privacy budgets ε and δ ; ; privacy-utility tradeoff coefficient θ ; sample size n ; local convex coefficient λ ; norm upper bound C ; cumulative Hessian upper bound H ; individual Hessian minimum eigenvalue upper bound ζ_{\min} ; bound looseness probability ρ .

$\tilde{\mathbf{w}}_0 \leftarrow \mathbf{w}^*$

$\mathcal{D}_{r_0} \leftarrow \mathcal{D}$

for $i = 1, \dots, k$ **do**

$\mathcal{D}_{r_i} \leftarrow \mathcal{D}_{r_{i-1}} \setminus \mathcal{D}_{f_i}$

$\mathbf{P}_{0,\lambda} \leftarrow \nabla \tilde{\mathbf{w}}_{i-1, \mathcal{K}, \mathcal{A}}$

for $t = 1, \dots, n$ **do**

 Sample X_{t_i} from \mathcal{D}_{f_i} with probability θ or sample X_{t_i} from \mathcal{D}_r with probability $1 - \theta$

if $X_{t_i} \sim \mathcal{D}_{f_i}$ **then**

$\mathbf{H}_{t,\lambda,i} \leftarrow \nabla_{\mathbf{w}}^2 \mathcal{L}_{\mathcal{K}}(\mathbf{w}^*, X_{t_i}) + \frac{\lambda I}{2\theta}$

else if $X_{t_i} \sim \mathcal{D}_r$ **then**

$\mathbf{H}_{t,\lambda} \leftarrow \nabla_{\mathbf{w}}^2 \mathcal{L}_{\mathcal{A}}(\mathbf{w}^*, X_{t_i}) + \frac{\lambda I}{2(1-\theta)}$

end if

$\mathbf{P}_{t,\lambda,i} = \mathbf{P}_{0,\lambda,i} + (\mathbf{I} - \frac{\mathbf{H}_{t,\lambda,i}}{H})\mathbf{P}_{t-1,\lambda,i}$.

end for

$\tilde{\mathbf{w}}_i \leftarrow \tilde{\mathbf{w}}_{i-1} - \frac{\mathbf{P}_{n,\lambda,i}}{H}$

end for

Compute Δ as the bound in Eq. (G.23).

$\sigma = \frac{\Delta}{k\varepsilon} \sqrt{2 \ln(1.25/\delta)}$

$\mathbf{w}^- \leftarrow \tilde{\mathbf{w}}_k + Y$ where $Y \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$

return \mathbf{w}^-

K EXPERIMENTAL DETAILS

K.1 DATASET DETAILS

MNIST: The MNIST dataset contains 70k 28x28 greyscale images of hand-drawn digits in 10 classes (Deng, 2012). We conduct our experiments with 49k training images and 21k test images. The classes are mutually exclusive.

Kuzushiji-MNIST: The Kuzushiji-MNIST (KMNIST) dataset contains 70k 28x28 greyscale images of Japanese kanji in 10 classes (Clanuwat et al., 2018). We conduct our experiments with 49k training images and 21k test images. The classes are mutually exclusive.

CIFAR10: The CIFAR10 dataset consists of 60k 32x32 color images in 10 classes. The classes are mutually exclusive and include airplanes, automobiles, birds, cats, deer, dogs, frogs, horses, ships, and trucks (Krizhevsky et al., 2009).

CIFAR-100: The CIFAR-100 dataset is similar to CIFAR-10 but contains 100 classes, each with 600 images, making a total of 60k 32x32 color images. The 100 classes are grouped into 20 superclasses, and each image comes with a “fine” label (the class to which it belongs) and a “coarse” label (the superclass to which it belongs) (Krizhevsky et al., 2009).

SVHN: The Street View House Numbers (SVHN) dataset (Netzer et al., 2011) contains images of double-digit numbers on house walls as colored 32x32 images. We load SVHN with 100 classes, corresponding to 10 * 10 for each digit. There are 73k training images, 26k testing images.

TinyImageNet: The TinyImageNet dataset (Le and Yang, 2015) contains 100000 64x64 images.

K.2 MODEL DETAILS

LogReg: A logistic regression model that has a single linear layer between inputs and outputs, followed by a softmax output function.

MLP: A two-layer ReLU feedforward neural network.

2214 **ResNet8:** A [1,1,1,0] residual network, , with standard convolutional blocks, as described in (He
2215 et al., 2016).

2216 **ResNet18:** A [2,2,2,2] residual network, with standard convolutional blocks, as described in (He
2217 et al., 2016).

2218 **ResNet50:** A [3, 4, 6, 3] residual network, with bottleneck convolutional blocks, as described in (He
2219 et al., 2016).

2221 **ViT_S_16:** A vision transformer with ≈ 20 million parameters, as detailed in (Dosovitskiy et al.,
2222 2021).

2223 **ViT_B_16:** A vision transformer with ≈ 80 million parameters, as detailed in (Dosovitskiy et al.,
2224 2021).

2226 K.3 BASELINE DETAILS

2227 We implement several baselines and provide the rationale for their use below:

2228 **Pretrained:** This is simply the pretrained model corresponding to whichever model and benchmark
2229 is specified. The rationale for using this is to demonstrate that we alter uniformity significantly from
2230 before without tarnishing accuracy for either the retain or test sets.

2231 **Retrained:** This is a model retrained over the retain set, performing exact unlearning. The rationale
2232 for using this is to demonstrate that our methods mimic unlearning in how we preserve accuracy, but
2233 induce uniformity in a way that unlearning does not.

2234 **Synthetic:** This method proceeds as follows: for each instance in the forget set, sample k instances
2235 from the ε -ball, with respect to the ℓ_2 norm, around that instance. Then, assign these k instances
2236 random labels from the label space, choosing labels uniformly at random. Do this for all forget
2237 set instances, yielding $|\mathcal{D}_f|k$ new instances. Then, append this to the retain set and retrain over
2238 this augmented dataset. This provides strong accuracy for simple baselines like MNIST while also
2239 inducing uniformity, providing an alternative, simple algorithm to compare our method against in
2240 terms of time elapsed.

2241 **Label Differential Privacy:** Specifically, we use the multi-stage training method of Ghazi et al.
2242 (2021) to obtain a model which is differentially private with respect to the labels—that is, an adversary
2243 cannot be sure whether the label they obtain is the true label. This is related to our work, albeit
2244 addresses a different threat model, as described in Appx. C. Still, we believe it is important to
2245 demonstrate that our method achieves privacy while not sacrificing utility to the extend that label
2246 differential privacy does, since it is a well-known method in the privacy literature that addresses
2247 a similar problem. For our experiments, we use the official repository with the hyperparameters
2248 reported in the paper: <https://github.com/google-research/label-dp>.

2251 K.4 HYPERPARAMETER DETAILS

2252 Please note that, throughout, we do not do extensive hyperparameter optimization, which may lead to
2253 improved performance.

2254 We use a standard train-test split of 70-30 throughout. Pretraining and synthetic training have the
2255 same hyperparameters as pretraining, unless mentioned otherwise. We use ADAM, with standard
2256 PyTorch hyperparameters aside from learning rate and weight decay, throughout. A batch size of
2257 128 is used for pretraining and also for the retain set in Alg. 1 throughout. **To ensure that gradients
2258 were well-aligned between the two conflicting objectives in Alg. 1, we implemented and ran gradient
2259 surgery (Yu et al., 2020) before each gradient step. Roughly, gradient surgery projects the gradient of
2260 one objective onto the normal plane of the other objective’s gradient only when the two gradients
2261 conflict, i.e. have negative dot product. The projected gradient is chosen uniformly at random
2262 between the two. This effectively removes the component of the gradient that would increase the
2263 other objective’s loss.** For all Alg. 1 experiments and Alg. 2 experiments, we use a forget set size of
2264 100 with a batch size (when loading the forget set into the finetuning in Alg. 1) of 10. We perform
2265 Alg. 1 for 100 epochs and finetune Alg. 2 for 50 epochs before running the certified Newton step. We
2266 generally use the forward KL divergence between model softmax outputs and the uniform distribution
2267 for $\mathcal{L}_{\mathcal{K}}$ and the cross entropy between model predictions and ground truth labels for $\mathcal{L}_{\mathcal{A}}$. For LogReg,

2268 we instead use the square loss between the uniform softmax probabilities and the model softmax
 2269 outputs, since the forward KL is not necessarily convex in w in this case, while the square loss is;
 2270 this allows us to use Alg. 2 with small λ . For our synthetic baseline, we use $\varepsilon = 8/255$ throughout,
 2271 where ε is the size of the ε -ball where we sample instances to assign random labels for retraining. For
 2272 the LabelDP baseline, we use the multi-stage training algorithm of Ghazi et al. (2021) throughout.
 2273
 2274 Early stopping is implemented by saving the model which first meets the early stopping conditions,
 2275 and continuing to see if any model performs better in terms of confidence distance while still meeting
 2276 the early stopping conditions specified below.
 2277
 2278 **Compute:** We use two RTX 6000 Ada Generation NVIDIA GPUs throughout. The most resource
 2279 intensive experiments are the LabelDP experiments, which take up most of the memory on both GPUs.
 2280 Besides those, the other experiments take up at most a fourth of the compute resources available on
 2281 one GPU. No experiments ran required more compute than these two GPUs provide.
 2282
 2283 **MNIST, LogReg Pretraining:** Epochs: 25. Learning rate: 0.01.
 2284
 2285 **MNIST, MLP Pretraining:** Epochs: 5. Learning rate: 0.01.
 2286
 2287 **MNIST, ResNet18 Pretraining:** Epochs: 2. Learning rate: 0.001.
 2288
 2289 **MNIST, LogReg Alg. 1:** Learning rate: 0.01
 2290
 2291 **MNIST, MLP Alg. 1:** Learning rate: 0.01
 2292
 2293 **MNIST, SVM Alg. 1:** Epochs: 100. Learning rate: 0.1. Optimized with SGD.
 2294
 2295 **MNIST, ResNet18 Alg. 1:** Learning rate: 0.001. Early stopping criterion of a confidence distance
 2296 < 0.32 and a retain accuracy of $> 90\%$.
 2297
 2298 **MNIST, LogReg Alg. 2:** $M = 1$. $C = 10$, pretrained with PGD with the same hyperparameters
 2299 as the standard pretraining. Since the losses are convex in w , $\lambda_{\min} = 0$. $\lambda = 0.0001$. Following
 2300 Zhang et al. (2024), we use the variance σ^2 as a hyperparameter, corresponding to a broad range of
 2301 choices of ε and δ . We choose $\sigma = 0.001$. This results in large ε when choosing small δ , as typical
 2302 in differential privacy (Dwork et al., 2014) and certified unlearning (Qiao et al., 2025). However, we
 2303 still observe good induced uniformity. **We also use a learning rate of $\alpha = 0.5$ in the Newton step,**
 2304 **which, following Appx. G and bounding the gradient with αF instead of F , results in roughly twice**
 2305 **as large of an ε .**
 2306
 2307 **MNIST, LogReg Alg. 3:** $\lambda = 1.5$, $J = 4.0$, $n = 1000$, $b = 100$, and otherwise same hyperparame-
 2308 ters as Alg. 2.
 2309
 2310 **MNIST, MLP Alg. 3:** $\lambda = 1.4$, $J = 2.5$, $n = 1000$, $b = 100$, and a learning rate of $\alpha = 0.35$.
 2311
 2312 **MNIST, SVM Alg. 3:** $\lambda = 0.5$, $J = 2.5$, $n = 5000$, $b = 200$, and a learning rate of $\alpha = 0.05$. We
 2313 also used a twice-differentiable surrogate for the hinge loss to replace the standard hinge loss, which
 2314 was used as the pretraining loss $\mathcal{L}_{\mathcal{A}}$.
 2315
 2316 **MNIST, LogReg Synthetic Baseline:** Sampled k instances for each forget set instance: 5.
 2317
 2318 **MNIST, ResNet18 Synthetic Baseline:** Sampled k instances for each forget set instance: 500.
 2319
 2320 **MNIST, LogReg LabelDP Baseline:** Epochs: 200. Batch size: 256. Random flip, random left-right
 2321 flip, and random cutout (8). SGD with learning rate 0.4 with momentum 0.9. $\varepsilon = 2.0$. Mixup for
 stage 1: 16. Mixup for stage 2: 8. Data split evenly between the two stages. Piecewise constant
 learning rate scheduler. These hyperparameters are chosen to match those in the best results of Ghazi
 et al. (2021). See Ghazi et al. (2021) for more details on these hyperparameters.
 MNIST, ResNet18 LabelDP Baseline: Same as the MNIST LogReg LabelDP hyperparameters,
 except with a weight decay of 0.0005 throughout.
 KMNIST, LogReg Pretraining: Epochs: 100. Learning rate: 0.01.
 KMNIST, MLP Pretraining: Epochs: 100. Learning rate: 0.001.
 KMNIST, ResNet18 Pretraining: Epochs: 12 Learning rate: 0.002.
 KMNIST, LogReg Alg. 1: Same as pretraining.

2322 **KMNIST, MLP Alg. 1:** Learning rate: 0.01.
 2323
 2324 **KMNIST, ResNet18 Alg. 1:** Learning rate: 0.002. Early stopping criterion of a confidence distance
 2325 < 0.32 and a retain accuracy of $> 99\%$.
 2326 **KMNIST, ResNet18 Synthetic Baseline:** Sampled k instances for each forget set instance: 500.
 2327 **KMNIST, ResNet18 LabelDP Baseline:** Same as the MNIST ResNet18 LabelDP hyperparameters.
 2328
 2329 **SVHN, ResNet50 Pretraining:** Epochs: 150. Learning rate: 0.001. Weight decay: 0.00005.
 2330 **SVHN, ResNet50 Alg. 1:** Same as pretraining.
 2331
 2332 **SVHN, ResNet50 Synthetic Baseline:** Sampled k instances for each forget set instance: 500.
 2333 **SVHN, ResNet50 LabelDP Baseline:** Same as MNIST ResNet18 LabelDP hyperparameters.
 2334
 2335 **CIFAR10, ResNet18 Pretraining:** Epochs: 200 with SGD with a momentum of 0.9. Learning rate:
 2336 0.1. Weight decay: 0.0005.
 2337 **CIFAR10, ResNet18 Alg. 1:** Same as pretraining. Early stopping criterion of a confidence distance
 2338 < 0.42 and a retain accuracy of $> 87\%$.
 2339 **CIFAR10, ResNet50 Pretraining:** Same as CIFAR10 ResNet18.
 2340
 2341 **CIFAR10, ResNet50 Alg. 1:** Same as pretraining. Early stopping criterion of a confidence distance
 2342 < 0.42 and a retain accuracy of $> 87\%$.
 2343 **CIFAR10, ResNet8 Pretraining:** Same as CIFAR10 ResNet18.
 2344
 2345 **CIFAR10, ResNet8 Alg. 1:** Same as CIFAR10 ResNet18.
 2346
 2347 **CIFAR10, ResNet18 Synthetic Baseline:** Sampled k instances for each forget set instance: 5000.
 2348
 2349 **CIFAR10, ResNet50 Synthetic Baseline:** Sampled k instances for each forget set instance: 5000.
 2350
 2351 **CIFAR10, ResNet18 LabelDP Baseline:** Same as MNIST ResNet18 LabelDP hyperparameters,
 2352 except with a batch size of 512.
 2353
 2354 **CIFAR10, ResNet50 LabelDP Baseline:** Same as CIFAR10 ResNet18 LabelDP.
 2355
 2356 **CIFAR10, ViT_S_16 Finetuning:** 8 epochs. Learning rate 0.0001 with AdamW.
 2357
 2358 **CIFAR10, ViT_B_16 Finetuning:** 10 epochs. Learning rate 0.0001 with AdamW.
 2359
 2360 **CIFAR100, ResNet50 Pretraining:** Same as CIFAR10 ResNet18.
 2361
 2362 **CIFAR100, ResNet50 Alg. 1:** Same as pretraining. Early stopping criterion of a confidence distance
 2363 < 0.42 and a retain accuracy of $> 87\%$.
 2364
 2365 **CIFAR100, ResNet8 Pretraining:** Same as CIFAR100 ResNet50.
 2366
 2367 **CIFAR100, ResNet8 Alg. 1:** Same as CIFAR100 ResNet50.
 2368
 2369 **CIFAR100, ResNet50 Synthetic Baseline:** Sampled instances: 5000.
 2370
 2371 **CIFAR100, ResNet50 LabelDP Baseline:** Same as CIFAR10 ResNet18 LabelDP. Please note that
 2372 our results differ from the results reported in the original paper of Ghazi et al. (2021); however, we
 2373 verified our results through several runs and used the official paper repository at <https://github.com/google-research/label-dp> with the hyperparameters reported in the paper.
 2374
 2375 **CIFAR100, ViT_S_16 Finetuning:** 30 epochs. Learning rate 0.002 with SGD with momentum 0.9.
 500 warmup steps with cosine scheduler.
 2376
 2377 **CIFAR100, ViT_S_16 Alg. 1:** 100 epochs. Learning rate 0.001 with SGD.
 2378
 2379 **CIFAR100, ViT_B_16 Finetuning:** 45 epochs. Learning rate 0.002 with SGD with momentum 0.9.
 500 warmup steps with cosine scheduler.
 2380
 2381 **CIFAR100, ViT_B_16 Alg. 1:** Same as CIFAR100 ViT_S_16.

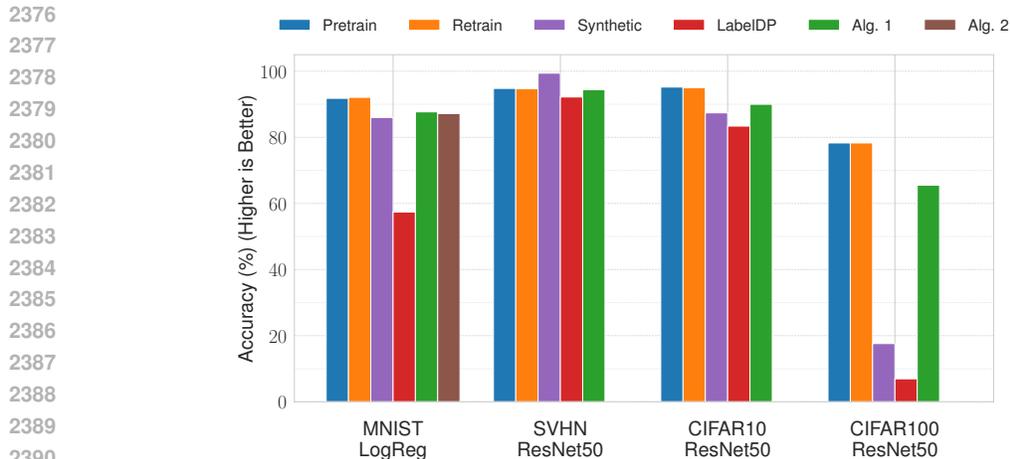


Figure L.5: Accuracy on test set for baselines as well as Alg. 1 and Alg. 2 with $\theta = 0.75$.

TinyImageNet, ViT_S_16 Finetuning: 30 epochs. Learning rate 0.0001, momentum 0.9, and weight decay 0.01 with SGD.

TinyImageNet, ViT_S_16 Alg. 1: Same as CIFAR100 ViT_S_16.

TinyImageNet, ViT_B_16 Finetuning: 50 epochs. Learning rate 0.0001, momentum 0.9, and weight decay 0.01 with SGD.

TinyImageNet, ViT_B_16 Alg. 1: Same as CIFAR100 ViT_S_16.

Attacks: $\alpha = 0.0001$. PGD learning rate: 0.001. PGD steps: 50.

Test-Set Finetuning: Finetune pretrained model for 20 more epochs with the same hyperparameters as pretraining, then run Alg. 1 with the same hyperparameters as the original run of Alg. 1. For CIFAR10/CIFAR100, finetune for 100 epochs.

L ADDITIONAL EXPERIMENTS

L.1 TEST SET ACCURACIES FOR MAIN PAPER EXPERIMENTS

We provide a figure, similar to Fig. 2a, for the test set in Fig. L.5. We observe similar results as one does on the retain set, with test accuracies preserved by Alg. 1 and Alg. 2.

L.2 TABLES FOR MAIN PAPER EXPERIMENTS

The tables for Fig. 2 are in Tab. L.1 and Tab. L.2. We include results for MNIST and KMNIST ResNet18 as well. We see that Alg. 1 induces uniformity without great damage to utility, while all other baselines—including the synthetic baseline—fail to do so without critically harming utility. Furthermore, we observe that ResNet50 performs better than ResNet18, providing more credibility to the claim in Sec. 6 that larger models tend to perform better when used in Alg. 1. Next, we observe that Alg. 1 can actually provide better retain and test accuracy than the pretrained model, as observed for ResNet50 over CIFAR100; this is because we also minimize the retain accuracy during finetuning. We similarly have to use early stopping for Alg. 1, as discussed in Sec. 6, since we use large models. Finally, we observe that LabelDP can induce uniformity, albeit at the cost of retain and test accuracy, but does so not only on the forget set but also the retain set; for a comparison of the confidence distances across the retain, test, and forget sets for LabelDP and our method, please see Tab. L.4. Additionally, for larger, more complex datasets like CIFAR100, LabelDP fails entirely. Please note that we do not perform extensive hyperparameter optimization during pretraining or retraining.

We observe similar results for Alg. 2 in Tab. L.2.

2430 Table L.1: Results for Alg. 1, used in Fig. 2. We find that we are able to induce uniformity while only
 2431 slightly decreasing retain and test accuracy. $\theta = 0.75$ throughout.

2432

2433

2434

2435

2436

2437

2438

2439

2440

2441

2442

2443

2444

2445

2446

2447

2448

2449

2450

2451

2452

2453

2454

2455

2456

2457

2458

2459

2460

2461

2462

2463

2464

2465

2466

2467

2468

2469

2470

2471

2472

2473

2474

2475

2476

2477

2478

2479

2480

2481

2482

2483

Dataset	Model	Method	Retain Acc.	Test Acc.	Conf. Dist. (Lower Better)
MNIST	ResNet18	Pretrain	98.0%	98.1%	0.877
		Retrain	97.3%	97.1%	0.876
		Synthetic	100.0%	99.1%	0.010
		LabelDP	98.8%	98.8%	0.593
		Alg. 1	99.6%	99.1%	0.070
KMNIST	ResNet18	Pretrain	98.2%	92.1%	0.880
		Retrain	98.4%	92.4%	0.884
		Synthetic	99.9%	96.7%	0.019
		LabelDP	98.9%	96.1%	0.530
		Alg. 1	99.1%	94.7%	0.257
SVHN	ResNet50	Pretrain	99.6%	94.8%	0.980
		Retrain	99.3%	94.7%	0.964
		Synthetic	99.9%	99.4%	0.013
		LabelDP	92.0%	92.2%	0.282
		Alg. 1	99.5%	94.4%	0.280
CIFAR10	ResNet18	Pretrain	100.0%	95.3%	0.898
		Retrain	100.0%	95.3%	0.891
		Synthetic	94.0%	89.7%	0.844
		LabelDP	85.8%	83.6%	0.359
		Alg. 1	89.6%	83.1%	0.377
	ResNet50	Pretrain	100.0%	95.2%	0.900
		Retrain	100.0%	95.0%	0.891
		Synthetic	91.4%	87.4%	0.818
		LabelDP	85.5%	83.4%	0.334
		Alg. 1	94.7%	90.0%	0.270
CIFAR100	ResNet50	Pretrain	100.0%	78.3%	0.902
		Retrain	100.0%	78.3%	0.765
		Synthetic	17.7%	17.6%	0.189
		LabelDP	8.41%	6.95%	0.203
		Alg. 1	91.4%	65.5%	0.298

2463

2464

2465 Table L.2: Results for Alg. 2 for logistic regression trained over MNIST, used in Fig. 2. $\theta = 0.75$
 2466 throughout.

2467

2468

2469

2470

2471

2472

2473

2474

2475

2476

2477

2478

2479

2480

2481

2482

2483

L.3 ADDITIONAL EXPERIMENTS ON TINYIMAGENET & ViT

We provide experimental results for Alg. 1 for ViT trained on CIFAR100 and TinyImageNet in Tab. L.3, observing similar behavior—in fact significantly lower confidence distance with little retain or test accuracy reduction—when compared to in Tab. L.1.

Table L.3: Results for Alg. 1 for ViT trained on CIFAR100 and TinyImageNet.

Dataset	Model	Method	Retain Acc.	Test Acc.	Conf. Dist. (Lower Better)
CIFAR100	ViT_S_16	Pretrain	95.2%	90.1%	0.942
		Retrain	93.4%	89.1%	0.883
		Synthetic Alg. 1	86.4%	84.8%	0.682
		Synthetic Alg. 1	91.6%	85.1%	0.036
	ViT_B_16	Pretrain	94.2%	91.2%	0.972
		Retrain	95.2%	91.0%	0.952
		Synthetic Alg. 1	17.7%	17.6%	0.189
		Synthetic Alg. 1	91.6%	88.6%	0.074
TinyImageNet	ViT_S_16	Pretrain	86.7%	84.4%	0.698
		Retrain	87.2%	84.4%	0.742
		Synthetic Alg. 1	87.9%	86.3%	0.833
		Synthetic Alg. 1	84.2%	81.4%	0.057
	ViT_B_16	Pretrain	95.0%	91.7%	0.822
		Retrain	96.8%	90.6%	0.826
		Synthetic Alg. 1	98.0%	84.4%	0.830
		Synthetic Alg. 1	91.8%	88.3%	0.037
		Pretrain	91.2%	83.5%	0.812
		Retrain	91.6%	80.9%	0.924
ResNet50	Synthetic Alg. 1	92.1%	80.6%	0.569	
	Synthetic Alg. 1	92.1%	81.6%	0.197	

Table L.4: A comparison of the confidence distances on the retain, test, and forget sets between Alg. 1 and LabelDP. In general, we induce uniformity on only the forget set, while maintaining confidently correct predictions on the retain and test sets, while LabelDP falls short. Note that this is only for one experiment run.

Dataset	Model	Method	Retain Conf. Dist. (Higher Better)	Test Conf. Dist. (Higher Better)	Forget Conf. Dist. (Lower Better)
MNIST	ResNet18	LabelDP	0.579	0.577	0.593
		Alg. 1	0.503	0.509	0.070
KMnist	ResNet18	LabelDP	0.495	0.466	0.530
		Alg. 1	0.870	0.828	0.257
SVHN	ResNet50	LabelDP	0.288	0.285	0.282
		Alg. 1	0.888	0.855	0.280
CIFAR10	ResNet18	LabelDP	0.371	0.365	0.359
		Alg. 1	0.724	0.690	0.377
	ResNet50	LabelDP	0.366	0.361	0.334
		Alg. 1	0.725	0.701	0.270
CIFAR100	ResNet50	LabelDP	0.182	0.156	0.203
		Alg. 1	0.576	0.470	0.298

L.4 LABELDP AND ALG. 1 CONFIDENCE DISTANCES FOR RETAIN, TEST, AND FORGET SETS

Here, we present Tab. L.4, which details the confidence distances for the retain, test, and forget sets of our method vs. LabelDP. Not only do we achieve better retain and test accuracy, but also we induce uniformity on *only* the forget set, while LabelDP induces uniformity on the forget, retain, and test sets altogether, functionally the same as adjusting the temperature. This does not suffice for our threat model, since we want to preserve confident predictions on the retain and test sets.

L.5 PARETO FRONTIER MAIN PAPER TABLE

The results which correspond to Fig. 3b, Fig. 3a, and Fig. L.12 are included in Tab. L.5 and Tab. L.6.

Table L.5: Results for Alg. 1 and Alg. 2 as we explore the Pareto frontier over MNIST, single run for Fig. 3b, Fig. 3a, and Fig. L.12.

Model	θ	Retain Acc.	Test Acc.	Conf. Dist. (Lower Better)
LogReg	0.000	93.7%	92.9%	0.817
	0.125	91.9%	91.2%	0.583
	0.250	92.9%	92.4%	0.178
	0.375	92.0%	92.3%	0.208
	0.500	92.4%	92.1%	0.342
	0.625	91.3%	91.1%	0.374
	0.750	91.6%	91.2%	0.263
	0.850	88.4%	87.5%	0.204
	0.950	89.2%	89.1%	0.169
Cert. LogReg	0.000	92.5%	92.2%	0.738
	0.125	92.2%	91.3%	0.573
	0.250	91.6%	91.5%	0.324
	0.375	91.4%	90.5%	0.327
	0.500	90.5%	90.6%	0.300
	0.625	90.6%	89.7%	0.451
	0.750	87.1%	87.2%	0.280
	0.850	89.3%	88.4%	0.206
	0.950	85.0%	85.7%	0.092
MLP	0.000	100.0%	98.1%	0.893
	0.125	98.3%	97.6%	0.399
	0.250	99.8%	97.4%	0.068
	0.375	97.4%	96.6%	0.247
	0.500	99.5%	97.1%	0.037
	0.625	96.4%	95.8%	0.166
	0.750	97.5%	95.7%	0.037
	0.850	92.5%	92.8%	0.096
	0.950	91.3%	90.0%	0.029
ResNet18	0.000	99.6%	99.4%	0.896
	0.125	99.3%	99.1%	0.892
	0.250	97.2%	97.5%	0.427
	0.375	99.6%	99.3%	0.586
	0.500	97.0%	97.0%	0.357
	0.625	97.0%	97.0%	0.357
	0.750	97.0%	97.0%	0.151
	0.850	95.7%	95.4%	0.234
	0.950	93.6%	94.0%	0.288

L.6 OPTIMIZATION DYNAMICS

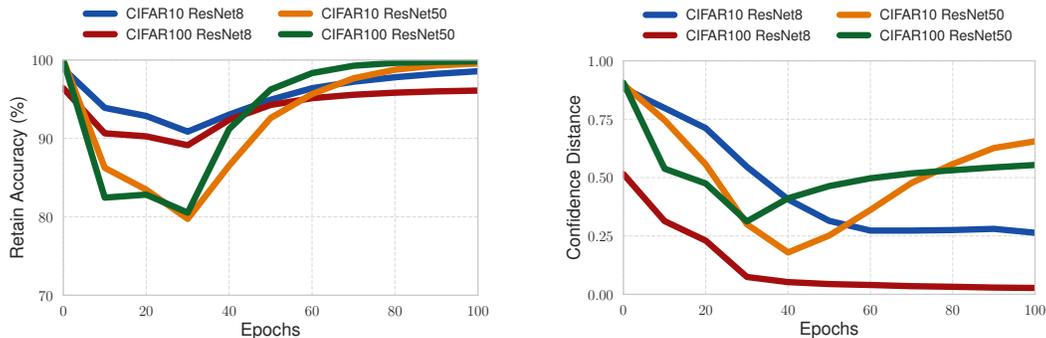
L.6.1 EMPIRICAL RESULTS

Upon using Alg. 1, we observe that logistic regression fails to induce uniformity for more complex benchmarks than MNIST, e.g. KMNIST. Logistic regression has poor test accuracy; we thus conclude that a model must be large enough to generalize well in order to have uniformity induced over it without a large cost to retain accuracy. We discuss mathematical intuition for this in Appx. L.6.2.

However, when using Alg. 1 on larger models, we observe that we need early stopping. After achieving a good uniformity-utility tradeoff, large models e.g. ResNet50 on CIFAR10 will powerfully increase accuracy at the cost of uniformity. This is undesired behavior when compared to, for example, a ResNet8 trained on CIFAR10, where we initially increase uniformity at the cost of accuracy but slowly regain accuracy without critically damaging uniformity. We illustrate this in Fig. L.6. Altogether, our method works best for large models with early stopping during finetuning. We characterize this mathematically in Appx. L.6.2 as well.

Table L.6: Results for Alg. 1 as we explore the Pareto frontier over CIFAR10 and CIFAR100 for ResNet50, single run for Fig. 3b, Fig. 3a, and Fig. L.12.

Dataset	θ	Retain Acc.	Test Acc.	Conf. Dist. (Lower Better)
CIFAR10	0.000	100.0%	92.9%	0.817
	0.125	99.9%	94.1%	0.616
	0.250	99.9%	93.8%	0.586
	0.375	99.9%	94.6%	0.501
	0.500	91.9%	85.7%	0.395
	0.625	90.9%	86.9%	0.343
	0.750	94.7%	90.0%	0.270
	0.850	56.8%	56.0%	0.108
	0.950	10.7%	10.4%	0.168
CIFAR100	0.000	92.5%	92.2%	0.738
	0.125	99.9%	77.0%	0.353
	0.250	99.9%	77.5%	0.252
	0.375	99.9%	77.0%	0.296
	0.500	99.9%	77.1%	0.301
	0.625	90.9%	70.2%	0.353
	0.750	91.4%	65.5%	0.298
	0.85	21.2%	20.6%	0.138
	0.950	40.8%	30.3%	0.019



(a) Retain Accuracy vs. Epochs, $\theta = 0.75$

(b) Confidence Distance vs. Epochs, $\theta = 0.75$

Figure L.6: For CIFAR10 and CIFAR100 ResNet50, we observe a sharp drop in confidence distance followed by a sharp increase in Fig. L.6b, in line with the drops and increases for retain accuracy in Fig. L.6a. Test accuracy is similar. This highlights the need for early stopping when using Alg. 1 for large models, since otherwise one escapes from a good privacy-utility tradeoff. For smaller models, e.g. MNIST MLP, this issue does not persist—we obtain good uniformity after an initial drop in accuracy, but then increase accuracy and decrease confidence distance simultaneously.

We provide a plot characterizing how test accuracy for Alg. 1 and Alg. 2 applied on CIFAR10 and CIFAR100 for $\theta = 0.75$ changes over 100 epochs in Fig. L.7 for ResNet8 and ResNet50. We observe similar behavior to retain accuracy.

Results as used in Fig. L.6a, Fig. L.6b, and Fig. L.7 are included in Tab. L.7.

L.6.2 INTUITION FOR EARLY STOPPING

In what follows, we give mathematical justification for the behavior observed in Fig. L.6 and Tab. L.7.

Firstly, recall that random vectors are nearly orthogonal in high dimensions (Vershynin, 2018). In particular, for larger models, the gradients will conflict i.e. point in opposite directions more strongly,

Table L.7: Studying the optimization dynamics of Alg. 1. Used in Fig. L.6.

Dataset	Model	Epoch	Retain Acc.	Test Acc.	Conf. Dist. (Lower Better)
CIFAR10	ResNet8	0	98.9%	90.8%	0.883
		10	88.9%	80.4%	0.713
		20	90.8%	83.1%	0.539
		30	92.9%	85.3%	0.386
		40	95.4%	87.3%	0.295
		50	96.5%	88.9%	0.263
		60	97.3%	89.9%	0.261
		70	97.9%	90.0%	0.295
		80	98.2%	90.4%	0.270
		90	98.6%	90.5%	0.276
100	98.9%	90.9%	0.245		
CIFAR10	ResNet50	0	100.0%	95.2%	0.899
		10	72.5%	66.4%	0.593
		20	77.9%	75.0%	0.176
		30	88.8%	85.3%	0.131
		40	92.9%	88.7%	0.231
		50	96.1%	91.1%	0.394
		60	98.0%	92.5%	0.459
		70	98.9%	93.2%	0.579
		80	99.4%	93.7%	0.636
		90	99.6%	94.0%	0.665
100	99.7%	94.0%	0.665		
CIFAR100	ResNet8	0	96.4%	67.8%	0.515
		10	84.9%	59.2%	0.112
		20	89.5%	63.4%	0.063
		30	93.0%	67.3%	0.048
		40	94.6%	68.5%	0.046
		50	95.2%	69.0%	0.038
		60	95.6%	69.4%	0.036
		70	95.9%	69.5%	0.031
		80	96.0%	69.8%	0.030
		90	96.1%	70.0%	0.026
100	96.2%	70.2%	0.026		
CIFAR100	ResNet50	0	100.0%	78.3%	0.906
		10	64.9%	50.8%	0.170
		20	83.6%	65.6%	0.348
		30	93.1%	71.7%	0.419
		40	96.9%	74.4%	0.467
		50	98.7%	75.5%	0.505
		60	99.4%	76.4%	0.519
		70	99.7%	76.7%	0.530
		80	99.8%	76.5%	0.546
		90	99.8%	76.8%	0.555
100	99.9%	77.3%	0.561		

2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753

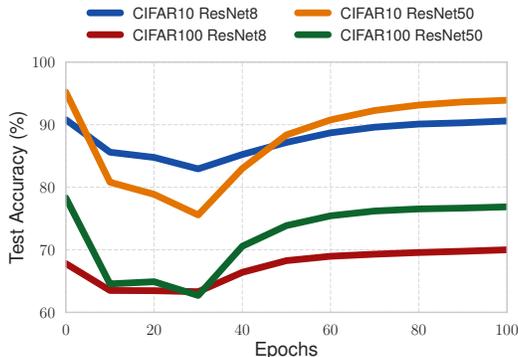


Figure L.7: Test Accuracy vs. Epochs, $\theta = 0.75$, MNIST. This has similar behavior to Fig. L.6a.

since their parameter space is very large. Second, every gradient step of our Alg. 1 is given by $\theta \nabla_{\mathbf{w}} \mathcal{L}_{\mathcal{K}}(\mathbf{w}, \mathcal{D}_f) + (1 - \theta) \nabla_{\mathbf{w}} \mathcal{L}_{\mathcal{A}}(\mathbf{w}, \mathcal{D}_r)$.

In what follows, consider a large model e.g. CIFAR10 ResNet50.

When we begin finetuning the pretrained model with Alg. 1, θ is large, $\mathcal{L}_{\mathcal{A}}$ is small, and $\mathcal{L}_{\mathcal{K}}$ is large. Thus, $\nabla_{\mathbf{w}} \mathcal{L}_{\mathcal{K}}$ significantly dominates $\nabla_{\mathbf{w}} \mathcal{L}_{\mathcal{A}}$. For large models, since the gradients are nearly orthogonal, we will move very fast in the direction of the forget gradient.

As finetuning continues, we reach a point where θ is large, $\mathcal{L}_{\mathcal{K}}$ is small, and $\mathcal{L}_{\mathcal{A}}$ is large. At this point, the $\nabla_{\mathbf{w}} \mathcal{L}_{\mathcal{A}}$ begins to dominate, albeit less significantly since θ is large. For large models, since the gradients are nearly orthogonal, we will move fast in the direction of the retain gradient. However, due to our choice of large θ , the retain gradient will be reduced in magnitude. Thus, we will move more slowly at this stage.

We must stop shortly after this, otherwise the forget loss will climb back up to a point where the model is no longer reasonably uniform.

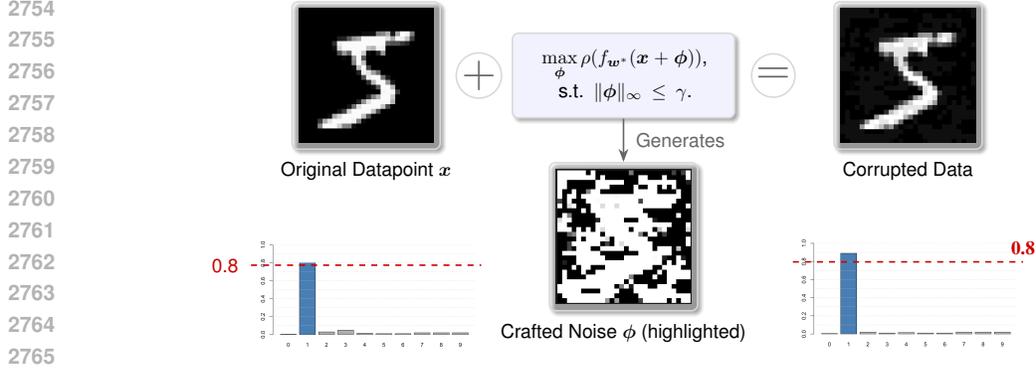
Importantly, since smaller models (e.g. CIFAR10 ResNet8) have smaller parameter spaces, the gradients do not conflict as much. Thus, when we begin finetuning the model, while we do increase in retain loss initially, after the uniform loss is minimized we can minimize the retain loss freely. As such, we can move in a direction that minimizes both the forget and retain gradients, and do not need to stop early.

However, as noted in the main paper, the model needs to be sufficiently large to achieve strong test accuracy, e.g. logistic regression trained on KMNIST does not work well. We hypothesize that this is because a model with more parameters has many more subspaces where two task losses can coexist in a way that provides a good tradeoff. We leave studying this and the above more formally to future work.

L.7 EVALUATING TEST-TIME PRIVACY ATTACKS

In what follows, we assume a test-time privacy (TTP) adversary with open-weight model access. We describe our attacks, specifically in Appx. A.

The results for the attacks with Alg. 5, Alg. 6, and Alg. 7 are in Tab. L.8 and Tab. L.9. We see that our method effectively defends against Alg. 5, Alg. 6, and Alg. 7 for various choices of γ in most scenarios. In our choices of γ , we follow the adversarial robustness literature, choosing γ sufficiently small such that the perturbation is invisible to the naked eye, but sufficiently large such that our attack is effective. However, in some cases, the attacks succeed despite the use of our method. Still, as demonstrated in Tab. L.10, we find that our algorithm renders the forget accuracy very low in these cases. As such, the adversary cannot be confidently correct—rather, in most cases, they can at best be confidently wrong. In particular, we have significantly better protection from attacks than in the pretrain, retrain, synthetic, or LabelDP cases. We provide visual intuition for our attacks in Fig. L.8.



2766 Figure L.8: We corrupt an instance to increase the confidence of the final prediction. Noise is
 2767 highlighted throughout for clarity.

2769 **Algorithm 5** Gaussian Noise Open-Weight Test Time Privacy Attack

2770 **Require:** Forget set \mathcal{D}_f ; pretrained model $w^* = \mathcal{A}(\mathcal{D})$; adversarial γ
 2771 $\mathcal{D}_f^{adv} = \emptyset$
 2772 **for** $i = 1, \dots, |\mathcal{D}_f|$ **do**
 2773 $x_{adv} = \mathcal{D}_f^{(i)} + \beta, \beta \sim \mathcal{N}(0, \gamma \mathbf{I})$
 2774 $\mathcal{D}_f^{adv, (i)} = x_{adv}$
 2775 **end for**
 2776 **return** \mathcal{D}_f^{adv}

2779
2780 L.8 ROBUSTNESS OF ALG. 1 CLASSIFIER ON NEIGHBORING TEST INSTANCES

2781
2782 To illustrate what happens for the finetuned classifier on neighboring test instances, we run an
 2783 experiment evaluating accuracy and average confidence distance on test instances which are nearest
 2784 neighbors to forget instances. Specifically, for each forget instance $x \in \mathcal{D}_f$, we found the nearest
 2785 neighbor of x in the test set. We evaluated this nearest neighbor in pixel space with respect to the ℓ_2
 2786 distance.

2787 Results are in Tab. L.11. We notice that the classifier works as intended. That is, we obtain high
 2788 accuracy as well as high confidence distance on the test set, including for nearby instances. While
 2789 we do observe a drop for CIFAR100 ResNet50, we also observe that the confidence distance is still
 2790 much higher than the 0.298 confidence distance on the forget set.

2791
2792 L.9 ENSURING TEST-TIME PRIVACY FOR TEST INSTANCES

2793
2794 In our paper, we focus on training data examples because this is the basis for scenarios addressed
 2795 by the GDPR, HIPAA, etc. Providing the same guarantee for non-training (test) data is an equally
 2796

2797 **Algorithm 6** FGSM-Style Open-Weight Test Time Privacy Attack

2798 **Require:** Forget set \mathcal{D}_f ; pretrained model $w^* = \mathcal{A}(\mathcal{D})$; adversarial γ ; symmetry breaking α
 2799 $\mathcal{D}_f^{adv} = \emptyset$
 2800 **for** $i = 1, \dots, |\mathcal{D}_f|$ **do**
 2801 $x_0 = \mathcal{D}_f^{(i)} + \beta, \beta \sim U([- \alpha \gamma, \alpha \gamma])$
 2802 $x_{adv} = \mathcal{D}_f^{(i)} + \gamma \text{sign}(\nabla_x \rho(f_{w^*}(x))) \Big|_{x=x_0}$
 2803 $\mathcal{D}_f^{adv, (i)} = x_{adv}$
 2804 **end for**
 2805 **return** \mathcal{D}_f^{adv}

Algorithm 7 PGD-Style Open-Weight Test Time Privacy Attack

```

2808 Require: Forget set  $\mathcal{D}_f$ ; pretrained model  $w^* = \mathcal{A}(\mathcal{D})$ ; adversarial  $\gamma$ ; symmetry breaking  $\alpha$ ; step
2809 count  $N$ 
2810  $\mathcal{D}_f^{adv} = \square$ 
2811 for  $i = 1, \dots, |\mathcal{D}_f|$  do
2812    $x_{adv}^0 = \mathcal{D}_f^{(i)} + \beta, \beta \sim U([- \alpha\gamma, \alpha\gamma])$ 
2813   for  $j = 1, \dots, N$  do
2814      $x_{adv}^j = x_{adv}^{j-1} + \beta \text{sign}(\nabla_x \rho(f_{w^*}(x)) \Big|_{x=x_0})$ 
2815      $x_{adv}^j = \Pi_{\mathcal{B}_\gamma(\mathcal{D}_f^{(i)})}(x_{adv}^j)$ 
2816   end for
2817    $\mathcal{D}_f^{adv, (i)} = x_{adv}^N$ 
2818 end for
2819 return  $\mathcal{D}_f^{adv}$ 

```

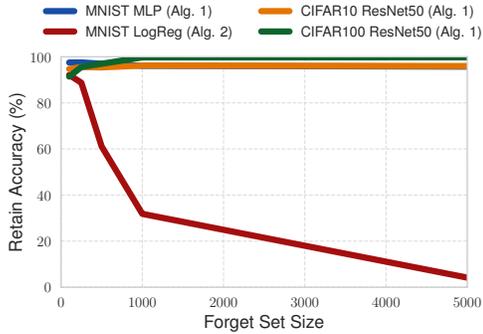
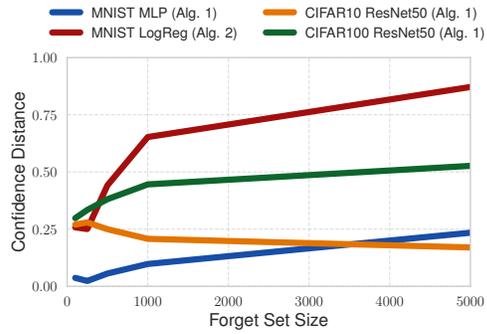
(a) Retain Accuracy vs. Epochs, $\theta = 0.75$ (b) Confidence Distance vs. Epochs, $\theta = 0.75$

Figure L.9: We observe that retain accuracy stays fairly stable as the forget size increases, in Fig. L.9a, except for Alg. 2 where it causes catastrophic failure due to the magnitude of the Newton step. Furthermore, we find that confidence distance slowly increases as the forget set size increases in Fig. L.9b.

important problem. However, the proposed method can be extended to cover this new case without loss of generality. One can just finetune on test instances highlighted to be corrupted and then run Alg. 1. In Tab. L.12, we find that finetuning with test instances yields similar performance to using our algorithm over just the training instances.

L.10 ABLATION STUDY ON FORGET SET SIZE

Figures are provided in Fig. L.9 and Fig. L.10. In Tab. L.13, we provide experiments for ResNet50 trained on CIFAR10 and CIFAR100 for Alg. 1. In Tab. L.14, we provide experiments on MLP trained over MNIST for Alg. 1 and logistic regression trained over MNIST for Alg. 2.

Throughout our experiments, we use a forget set size of 100. We do so because for our use case, it is likely that a data controller would want to induce uniformity only for a small number of instances. We observe that as one increases the forget set size, it becomes harder to induce uniformity with the same hyperparameters. Still, for Alg. 1, we are able to obtain strong uniformity with good retain and test accuracy for significantly larger forget set sizes. Furthermore, we observe that Alg. 2 fails for sufficiently large forget set size; this is likely because Hessian matrix is significantly larger (in norm) for a larger forget set, resulting in a catastrophically large Newton step. Mitigating this phenomenon is left to future work, where Hessian-free techniques like those of Qiao et al. (2025) may be advantageous.

2862 Table L.8: Confidence distances over the forget set after Alg. 5, Alg. 6, and Alg. 7 are applied to
 2863 pretrained models and models finetuned with Alg. 1 for MNIST and KMNIST. Lower is better.
 2864

2865	2866	2867	2868	2869	2870	2871	2872	2873	2874	2875	2876	2877	2878	2879	2880	2881	2882	2883	2884	2885	2886	2887	2888	2889	2890	2891	2892	2893	2894	2895	2896	2897	2898	2899	2900
Dataset	Model	γ	Method	Attack	Prior Conf. Dist. (Lower Better)	Attack Conf. Dist. (Lower Better)																													
MNIST	MLP	$\frac{2}{255}$	Pretrain	Alg. 6	0.879	0.884																													
			Alg. 5	0.037	0.043																														
			Alg. 1	Alg. 6	0.037	0.054																													
			Alg. 7	0.037	0.185																														
			Pretrain	Alg. 6	0.879	0.888																													
			Alg. 5	0.037	0.064																														
		$\frac{5}{255}$	Alg. 1	Alg. 6	0.037	0.089																													
			Alg. 7	0.037	0.488																														
			Pretrain	Alg. 6	0.879	0.891																													
			Alg. 5	0.037	0.091																														
			$\frac{8}{255}$	Alg. 1	Alg. 6	0.037	0.125																												
				Alg. 7	0.037	0.632																													
	ResNet18	$\frac{2}{255}$	Pretrain	Alg. 6	0.895	0.896																													
			Alg. 5	0.070	0.075																														
			Alg. 1	Alg. 6	0.070	0.133																													
			Alg. 7	0.070	0.133																														
			Pretrain	Alg. 6	0.895	0.897																													
			Alg. 5	0.070	0.088																														
		$\frac{5}{255}$	Alg. 1	Alg. 6	0.070	0.164																													
			Alg. 7	0.070	0.350																														
			Pretrain	Alg. 6	0.895	0.898																													
			Alg. 5	0.070	0.116																														
			$\frac{8}{255}$	Alg. 1	Alg. 6	0.070	0.248																												
				Alg. 7	0.070	0.638																													
KMNIST	ResNet18	$\frac{2}{255}$	Pretrain	Alg. 6	0.858	0.865																													
			Alg. 5	0.257	0.258																														
			Alg. 1	Alg. 6	0.257	0.302																													
			Alg. 7	0.257	0.317																														
			Pretrain	Alg. 6	0.858	0.872																													
			Alg. 5	0.257	0.259																														
		$\frac{5}{255}$	Alg. 1	Alg. 6	0.257	0.370																													
			Alg. 7	0.257	0.420																														
			Pretrain	Alg. 6	0.858	0.878																													
			Alg. 5	0.257	0.259																														
			$\frac{8}{255}$	Alg. 1	Alg. 6	0.257	0.433																												
				Alg. 7	0.257	0.520																													

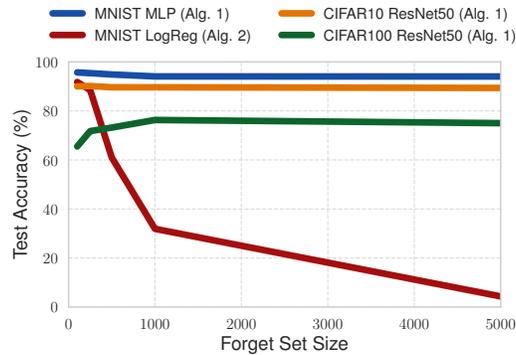
2901
 2902
 2903 The results in figures Fig. L.9 and Fig. L.10, as well as Tab. L.14 and Tab. L.13, we keep the batch size
 2904 used for the forget set in Alg. 1 to 10 instances. However, upon scaling the batch size proportionately
 2905 with the forget set size (specifically dividing the forget set size by ten), we observe much better
 2906 performance. These results are in Tab. L.15. Specifically, for a MLP trained on MNIST and a forget
 2907 set size of 5000, we obtain much better uniformity than when we did not scale the batch size.
 2908

2909 L.11 EVALUATING CONFIDENCE DISTANCE AS A TTP METRIC

2910
 2911 The ℓ_2 metric $\|f(\mathbf{x}) - \frac{\vec{1}}{K}\|_2$ has similar utility to our presented metric. However, it is slightly less
 2912 interpretable and may accidentally overpenalize uncertain outputs. For example, if one class has
 2913 no probability but the other 9 classes are uniform. Still, as demonstrated in Tab. L.16, we find that
 2914 we minimize this metric as well for the same models and datasets. This holds similarly for other
 2915 potential metrics e.g. the ℓ_1 metric.

2916 Table L.9: Confidence distances over the forget set after Alg. 5, Alg. 6, and Alg. 7 are applied to
 2917 pretrained models and models finetuned with Alg. 1 for SVHN, CIFAR10, and CIFAR100. Lower is
 2918 better.
 2919

2920 Dataset	2921 Model	γ	Method	Attack	Prior Conf. Dist. (Lower Better)	Attack Conf. Dist. (Lower Better)
2922 SVHN	ResNet50	$\frac{1}{255}$	Pretrain	Alg. 6	0.972	0.886
			Alg. 1	Alg. 5	0.289	0.519
			Alg. 1	Alg. 6	0.289	0.571
		$\frac{2}{255}$	Pretrain	Alg. 6	0.289	0.582
			Alg. 1	Alg. 5	0.972	0.904
			Alg. 1	Alg. 6	0.289	0.519
2927 CIFAR10	ResNet18	$\frac{1}{255}$	Pretrain	Alg. 6	0.289	0.613
			Alg. 1	Alg. 6	0.289	0.640
			Alg. 1	Alg. 7	0.289	0.640
		$\frac{2}{255}$	Pretrain	Alg. 6	0.898	0.898
			Alg. 1	Alg. 5	0.377	0.300
			Alg. 1	Alg. 6	0.377	0.353
	ResNet50	$\frac{1}{255}$	Pretrain	Alg. 6	0.377	0.356
			Alg. 1	Alg. 6	0.898	0.822
			Alg. 1	Alg. 7	0.377	0.301
		$\frac{2}{255}$	Pretrain	Alg. 6	0.377	0.397
			Alg. 1	Alg. 6	0.377	0.397
			Alg. 1	Alg. 7	0.377	0.408
2937 CIFAR100	ResNet50	$\frac{1}{255}$	Pretrain	Alg. 6	0.900	0.806
			Alg. 1	Alg. 5	0.270	0.336
			Alg. 1	Alg. 6	0.270	0.374
		$\frac{2}{255}$	Pretrain	Alg. 6	0.270	0.378
			Alg. 1	Alg. 6	0.900	0.827
			Alg. 1	Alg. 7	0.270	0.336
	ResNet50	$\frac{1}{255}$	Pretrain	Alg. 6	0.270	0.407
			Alg. 1	Alg. 6	0.270	0.425
			Alg. 1	Alg. 7	0.270	0.425
		$\frac{2}{255}$	Pretrain	Alg. 6	0.906	0.587
			Alg. 1	Alg. 5	0.298	0.275
			Alg. 1	Alg. 6	0.298	0.360
ResNet50	$\frac{1}{255}$	Pretrain	Alg. 6	0.298	0.373	
		Alg. 1	Alg. 6	0.298	0.373	
		Alg. 1	Alg. 7	0.298	0.373	
	$\frac{2}{255}$	Pretrain	Alg. 6	0.906	0.652	
		Alg. 1	Alg. 5	0.298	0.276	
		Alg. 1	Alg. 6	0.298	0.421	
ResNet50	$\frac{2}{255}$	Pretrain	Alg. 6	0.298	0.468	
		Alg. 1	Alg. 6	0.298	0.468	
		Alg. 1	Alg. 7	0.298	0.468	



2966 Figure L.10: Test Accuracy vs. Forget Set Size, $\theta = 0.75$. This has similar behavior to Fig. L.9a.
 2967
 2968
 2969

2970 Table L.10: Accuracies over the forget set after Alg. 6 and Alg. 7 are applied to pretrained models and
 2971 models finetuned with Alg. 1 for SVHN, CIFAR10, and CIFAR100. We see that Alg. 1 significantly
 2972 lowers forget set accuracy.
 2973

2974	Dataset	Model	γ	Method	Attack	Prior. Forget Acc. (Lower Better)	Atk. Forget Acc. (Lower Better)	
2976	MNIST	MLP	$\frac{5}{255}$	Pretrain	Alg. 7	97.0%	96.0%	
2977				Alg. 1	Alg. 7	71.0%	43.0%	
2978			$\frac{8}{255}$	Pretrain	Alg. 7	98.3%	96.0%	
2979				Alg. 1	Alg. 7	71.0%	41.0%	
2980		ResNet18	$\frac{8}{255}$	Pretrain	Alg. 7	98.9%	100.0%	
2981				Alg. 1	Alg. 7	28.0%	54.0%	
2982	KMNIST	ResNet18	$\frac{8}{255}$	Pretrain	Alg. 7	96.0%	96.0%	
2983				Alg. 1	Alg. 7	50.0%	55.0%	
2984	SVHN	ResNet50	$\frac{1}{255}$	Pretrain	Alg. 6	97.0%	66.0%	
2985				Alg. 5	40.0%	53.0%		
2986				Alg. 1	Alg. 6	40.0%	52.0%	
2987				Alg. 7	40.0%	52.0%		
2988			$\frac{2}{255}$	Pretrain	Alg. 6	97.0%	66.0%	
2989				Alg. 5	40.0%	52.0%		
2990				Alg. 1	Alg. 6	40.0%	53.0%	
2991				Alg. 7	40.0%	53.0%		
2992	CIFAR10	ResNet18	$\frac{1}{255}$	Pretrain	Alg. 6	100.0%	61.0%	
2993				Alg. 5	60.0%	35.0%		
2994				Alg. 1	Alg. 6	60.0%	34.0%	
2995			$\frac{2}{255}$	Pretrain	Alg. 6	100.0%	61.0%	
2996				Alg. 5	60.0%	35.0%		
2997				Alg. 1	Alg. 6	60.0%	33.0%	
2998		ResNet50	$\frac{1}{255}$	Pretrain	Alg. 6	100.0%	56.0%	
2999				Alg. 5	55.0%	33.0%		
3000				Alg. 1	Alg. 6	55.0%	30.0%	
3001			$\frac{2}{255}$	Pretrain	Alg. 6	100.0%	56.0%	
3002				Alg. 5	55.0%	30.0%		
3003				Alg. 1	Alg. 6	55.0%	30.0%	
3004	CIFAR100	ResNet50	$\frac{1}{255}$	Pretrain	Alg. 6	100.0%	32.0%	
3005				Alg. 5	48.0%	11.0%		
3006				Alg. 1	Alg. 6	48.0%	10.0%	
3007		$\frac{2}{255}$	Pretrain	Alg. 6	100.0%	32.0%		
3008			Alg. 5	48.0%	11.0%			
3009			Alg. 1	Alg. 6	48.0%	10.0%		
3010				Alg. 7	48.0%	11.0%		
3011							48.0%	10.0%
3012							48.0%	11.0%

3015 L.12 AN ADDITIONAL BASELINE WITH RANDOMLY SAMPLED LABELS: GAUSSIANUNIFORM

3016 In what follows, we present the *GaussianUniform* baseline, an alternative idea to our approach based
 3017 on the notable work of Zhang et al. (2017), which demonstrates that a neural network can fully
 3018 minimize its loss over a training dataset where samples have labels sampled uniformly at random.
 3019 The approach of GaussianUniform is as follows:
 3020

- 3021 1. Begin with a training dataset $\mathcal{D} = \mathcal{D}_f \cup \mathcal{D}_r$.
- 3022 2. Perturb all samples in \mathcal{D} to yield $\mathcal{D}' = \mathcal{D}'_f \cup \mathcal{D}'_r$. We use mean zero Gaussian noise with 0.1
 3023 variance, which adds a small amount of noise.

Table L.11: Accuracies and confidence distances for test instances which are nearest neighbors (with respect to ℓ_2 distance) of forget set instances. We observe that models continue to confidently and correctly classify nearby test instances after finetuning with Alg. 1.

Dataset	Model	Method	Acc.	Conf. Dist. (Higher Better)
MNIST	MLP	Pretrain	100.0%	0.894
		Alg. 1	93.0%	0.754
	ResNet18	Pretrain	100.0%	0.896
		Alg. 1	100.0%	0.875
KMNIST	ResNet18	Pretrain	99.0%	0.871
		Alg. 1	99.0%	0.850
SVHN	ResNet50	Pretrain	95.0%	0.982
		Alg. 1	95.0%	0.939
CIFAR10	ResNet18	Pretrain	98.0%	0.876
		Alg. 1	85.0%	0.538
	ResNet50	Pretrain	96.0%	0.884
		Alg. 1	90.0%	0.700
CIFAR100	ResNet50	Pretrain	78.0%	0.778
		Alg. 1	66.0%	0.484

3. Sample all labels in \mathcal{D} uniformly at random to yield $\tilde{\mathcal{D}} = \tilde{\mathcal{D}}_f \cup \tilde{\mathcal{D}}_r$.

4. Train $\mathcal{A}(\mathcal{D}' \cup \tilde{\mathcal{D}})$.

In this scenario, inducing uncertainty may not be necessary, since the forget set would have very strong uniformity, with strong accuracy on the retain set available by slightly perturbing with Gaussian noise. We use the same hyperparameters as pretraining for the respective model and dataset tested, as reported in Appx. K. We find that this is not the case for ResNet50 trained on SVHN, CIFAR10, and CIFAR100. However, it achieves very poor test accuracy compared to Alg. 1 on both normal $\mathcal{D}_{\text{test}}$ and perturbed test $\mathcal{D}'_{\text{test}}$ datasets; thus, we prefer Alg. 1 to this approach, since it generalizes well and also does not require retraining. Results are in Tab. L.17.

L.13 TIGHTNESS OF BOUND IN THEOREM 5.5

To evaluate how tight our bound is, we run an experiment for MNIST logistic regression. We use the notation of theorem 5.5 in Tab. L.18. We find that our constant bound is fairly tight as $\theta \rightarrow 1$; we leave using more advanced techniques to ensure better tightness to future work.

L.14 CONFIDENCE INTERVALS FOR MAIN PAPER EXPERIMENTS

In what follows, we report confidence intervals for only ResNet50 trained on SVHN due to compute constraints. We find that variance is low in Tab. L.19.

L.15 PROPORTIONS OF TIME ELAPSED IN ALG. 1

Results are reported in Tab. L.20.

L.16 WARMUP VALUES FOR MNIST LOGREG

Results are contained in Tab. L.21. We find that after applying the certified Newton step in Alg. 2, we obtain better retain and test accuracy, at small cost to uniformity. Thus, warming up is not the only component of achieving good results in Alg. 2.

3132 Table L.13: Results on applying Alg. 1 on ResNet50 for various forget set sizes over CIFAR10 and
 3133 CIFAR100. Please see Appx. L.10 for a discussion on Alg. 2. $\theta = 0.75$ throughout.
 3134

3135

3136 Dataset	3136 Model	3136 Forget Size	3136 Retain Acc.	3136 Test Acc.	3136 Conf. Dist. (Lower Better)
3137 CIFAR10	3137 ResNet50	3137 100	3137 94.7%	3137 90.0%	3137 0.270
		3138 250	3138 96.4%	3138 90.3%	3138 0.289
		3139 500	3139 95.1%	3139 88.6%	3139 0.190
		3140 1000	3140 97.0%	3140 90.0%	3140 0.144
		3141 5000	3141 95.8%	3141 89.5%	3141 0.176
3142 CIFAR100	3142 ResNet50	3142 100	3142 91.4%	3142 65.5%	3142 0.298
		3143 250	3143 99.7%	3143 78.0%	3143 0.370
		3144 500	3144 99.8%	3144 76.1%	3144 0.475
		3145 1000	3145 99.7%	3145 74.8%	3145 0.492
		3146 5000	3146 99.8%	3146 74.1%	3146 0.613

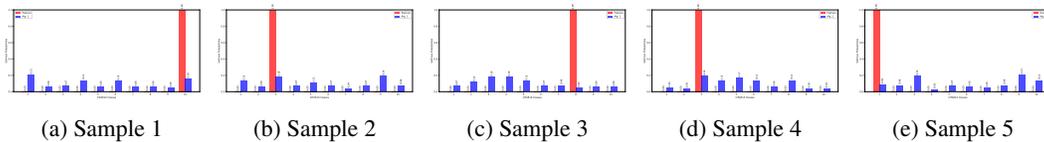
3147

3148 Table L.14: Results on applying Alg. 1 and Alg. 2 on various forget set sizes over MNIST. We observe
 3149 that, while Alg. 1 still works well, confidence distance increases as forget set size does; please see
 3150 Appx. L.10 for a discussion on Alg. 2. $\theta = 0.75$ throughout.
 3151

3152

3152 Method	3152 Model	3152 Forget Size	3152 Retain Acc.	3152 Test Acc.	3152 Conf. Dist. (Lower Better)
3153 Alg. 1	3153 MLP	3153 100	3153 97.5%	3153 95.7%	3153 0.037
		3154 250	3154 97.6%	3154 95.1%	3154 0.010
		3155 500	3155 95.8%	3155 93.8%	3155 0.121
		3156 1000	3156 94.8%	3156 93.3%	3156 0.162
		3157 5000	3157 96.6%	3157 95.0%	3157 0.420
3158 Alg. 2	3158 LogReg	3158 100	3158 92.1%	3158 91.8%	3158 0.258
		3159 250	3159 85.4%	3159 85.2%	3159 0.243
		3160 500	3160 5.7%	3160 5.9%	3160 0.824
		3161 1000	3161 4.4%	3161 4.6%	3161 0.891
		3162 5000	3162 2.4%	3162 2.5%	3162 0.899

3163



3169

3170 Figure L.11: Comparison of Pretrain (red) and Alg. 1 (blue) softmax probabilities across five different
 3171 CIFAR10 forget set samples.
 3172

3173

3174

3175

3176

3177

3178

3179

3180

3181

3182

3183

3184

3185

3176 Table L.15: Results on applying Alg. 1 on various forget set sizes over MNIST while scaling the
 3177 batches. We observe much better performance for a forget set size of 5000 than in Tab. L.14, and
 3178 comparable performance on other forget set sizes.
 3179

3180 Method	3180 Model	3180 Forget Size	3180 Retain Acc.	3180 Test Acc.	3180 Conf. Dist. (Lower Better)
3181 Alg. 1	3181 MLP	3181 100	3181 97.5%	3181 95.7%	3181 0.037
		3182 250	3182 95.6%	3182 96.6%	3182 0.131
		3183 500	3183 93.8%	3183 93.4%	3183 0.187
		3184 1000	3184 92.3%	3184 92.2%	3184 0.192
		3185 5000	3185 95.3%	3185 94.9%	3185 0.227

Table L.16: ℓ_2 confidence distances for models finetuned with Alg. 1.

Dataset	Model	Conf. Dist. Type	Pretrained (Lower Better)	Alg. 1 (Lower Better)
MNIST	MLP	Paper	0.879	0.037
		ℓ_2	0.930	0.053
MNIST	ResNet18	Paper	0.895	0.070
		ℓ_2	0.944	0.070
KMNIST	ResNet18	Paper	0.880	0.257
		ℓ_2	0.911	0.302
SVHN	ResNet50	Paper	0.972	0.289
		ℓ_2	0.979	0.298
CIFAR10	ResNet18	Paper	0.898	0.377
		ℓ_2	0.947	0.435
CIFAR10	ResNet50	Paper	0.900	0.270
		ℓ_2	0.948	0.323
CIFAR100	ResNet50	Paper	0.902	0.298
		ℓ_2	0.911	0.311

Table L.17: Results for the *GaussianUniform* baseline described in Appx. L.12. This method results in significantly degraded accuracy on the test set compared to a model finetuned with Alg. 1, despite achieving high accuracy on the perturbed retain set (\mathcal{D}'_r). Note that performance on $\tilde{\mathcal{D}}_r$ is similar to \mathcal{D}'_r .

Dataset	Model	Train Set Acc. (%)		Test Set Acc. (%)			Conf. Dist.
		\mathcal{D}'_r	$\tilde{\mathcal{D}}'_f$	$\mathcal{D}_{\text{test}}$	$\tilde{\mathcal{D}}'_{\text{test}}$	Alg. 1	\mathcal{D}_f
SVHN	ResNet50	82.4%	2.0%	6.2%	81.1%	94.4%	0.009
CIFAR10	ResNet50	100.0%	12.0%	11.1%	73.3%	90.0%	0.573
CIFAR100	ResNet50	99.9%	4.0%	3.4%	40.4%	65.5%	0.057

L.18 RESULTS FOR KMNIST LOGREG AND MLP

Results are contained in Tab. L.22. As mentioned in Sec. 6, one can see that a small model for a more complex benchmark, logistic regression on KMNIST, fails to induce uniformity, since the pretrained model is too small to generalize. However, on a bigger model i.e. an MLP trained over KMNIST, since it is large enough to generalize, one can induce uniformity over it. Thus, larger models which generalize well are preferred for our method, in line with the goals of ML.

L.19 ABLATION STUDY ON SYNTHETIC BASELINE SAMPLE SIZE

Below, we study what we happen if we increase the number of samples sampled in the ε -ball in the synthetic baseline. For each forget set instance, we sample k instances from the ε -ball around the forget set, and then assign random labels to these instances, yielding an additional $|\mathcal{D}_f|k$ instances in the training data. We then retrain the model over the retain set along with these new $|\mathcal{D}_f|k$ instances. For a MLP trained over MNIST, we observe better performance as we increase sample size. However, for a ResNet18 trained over CIFAR10, even if we have a very large sample size. This is presented in Tab. L.23. Thus, since Alg. 1 can induce uniformity as shown in Tab. L.1, without great cost to retain or test accuracy, it is better than the synthetic baseline.

Table L.18: Comparison of bounds

$\alpha^* - \alpha(1)$	O(1) bound size
1.678	3.374

3240 Table L.19: Results for Alg. 1 for ResNet50 trained on SVHN over three runs. We find that we are
 3241 able to induce uniformity while only slightly decreasing retain and test accuracy, $\theta = 0.75$ throughout,
 3242 with minimal variance for all metrics.

Dataset	Model	Method	Retain Acc.	Test Acc.	Conf. Dist. (Lower Better)
SVHN	ResNet50	Pretrain	99.9% \pm 0.1%	95.3% \pm 0.5%	0.987 \pm 0.001
		Alg. 1	99.5% \pm 0.2%	94.8% \pm 0.3%	0.276 \pm 0.004

3248 Table L.20: Time proportions of each step in Alg. 1.

Dataset	Model	Retain Grad.	Forget Grad.	Surgery	Reg. Grad.	Step
MNIST	MLP	0.967	0.033	0.000	0.000	0.000
	ResNet18	0.984	0.016	0.0010	0.000	0.000
KMNIST	ResNet18	0.989	0.011	0.000	0.000	0.000
SVHN	ResNet50	0.980	0.020	0.000	0.000	0.000
CIFAR10	ResNet18	0.989	0.011	0.000	0.000	0.000
	ResNet50	0.992	0.008	0.000	0.000	0.000
CIFAR100	ResNet50	0.993	0.006	0.000	0.000	0.000

3251
3252
3253
3254
3255
3256
3257
3258
3259
3260 L.20 TEST ACCURACY PLOT FOR PARETO FRONTIER EXPERIMENTS

3261 We provide a plot characterizing test accuracy for Alg. 1 and Alg. 2 applied on MNIST for various
 3262 choices of θ in Fig. L.12.

3263
3264
3265 L.21 RESULTS FOR ALG. 3

3266 Due to the issues presented in Appx. F, we focus on Alg. 1 and Alg. 2 in the main paper and present
 3267 results for Alg. 3 here. In Tab. L.24, we show that Alg. 1 can induce uniformity while maintaining
 3268 accuracy for logistic regression, albeit with a stronger tradeoff than Alg. 2

3269
3270
3271 L.22 HYPERPARAMETER SENSITIVITY FOR ALG. 3

3272 In what follows, we demonstrate the sensitivity of Alg. 3 to the choice of hyperparameters for
 3273 logistic regression trained on MNIST (using the forward KL for uniformity). Here $b = 1$ and
 3274 $n = 1000$. Below, the x -axis is the estimator recursive timesteps $t = 1, \dots, n$ and the y axis is the
 3275 difference between the exact Newton step and the estimated Newton step in Frobenius norm (before
 3276 any Gaussian noise is added). Specifically, we find that if $\lambda \approx |\lambda_{\min}|$ and $J \approx |\lambda + \lambda_{\min}|$ does not
 3277 hold, the estimator diverges very quickly. Specifically, as shown in Fig. L.13, if $\lambda = 1.0, J = 2.5$
 3278

3279
3280 Table L.21: Warmup values for Alg. 2 for logistic regression trained over MNIST, contrasted with the
 3281 values after Alg. 2 is applied.

θ	Method	Retain Acc.	Test Acc.	Conf. Dist. (Lower Better)
0.0	Warmup	91.4%	91.6%	0.765
0.0	Alg. 2	92.5%	92.2%	0.738
0.25	Warmup	90.1%	90.4%	0.283
0.25	Alg. 2	91.6%	91.5%	0.324
0.50	Warmup	89.3%	89.7%	0.215
0.50	Alg. 2	90.5%	90.6%	0.300
0.75	Warmup	88.4%	88.6%	0.154
0.75	Alg. 2	87.1%	87.2%	0.280
0.95	Warmup	85.4%	86.0%	0.097
0.95	Alg. 2	85.0%	85.7%	0.092

3294 Table L.22: Results for Alg. 1 applied to logistic regression and MLP trained over KMNIST, $\theta = 0.75$
 3295 for Alg. 1.
 3296

3297

Model	Method	Retain Acc.	Test Acc.	Conf. Dist. (Lower Better)
LogReg	Pretrain	81.4%	66.4%	0.775
	Retrain	80.9%	65.3%	0.770
	Alg. 1	77.4%	63.4%	0.770
MLP	Pretrain	100%	88.4%	0.900
	Retrain	100%	88.5%	0.887
	Alg. 1	92.8%	80.3%	0.039

3305
3306

3307 Table L.23: Results for the synthetic baseline applied with various sampled k on a MLP over
 3308 MNIST and a ResNet18 over CIFAR10. We observe that increasing k yields better performance,
 3309 but nevertheless even very large k (an additional 50k instances, with a forget set size of 100) fails to
 3310 induce uniformity for CIFAR10.
 3311

3312

Dataset	Model	Sampled k	Retain Acc.	Test Acc.	Conf. Dist. (Lower Better)
MNIST	MLP	5	99.4%	97.1%	0.541
		25	99.0%	96.4%	0.183
		125	99.4%	96.8%	0.105
		250	98.6%	96.0%	0.066
		500	99.6%	96.3%	0.003
CIFAR10	ResNet18	5	99.0%	91.0%	0.683
		25	99.2%	91.3%	0.865
		125	98.8%	90.9%	0.856
		250	98.4%	90.7%	0.869
		500	98.3%	91.1%	0.852
		5000	94.0%	89.7%	0.844

3324
3325
3326

3327 (here, λ is off by 0.5) or $\lambda = 0.5, J = 2.0$ (here, J is off by 0.5), the estimator diverges before
 3328 timestep 150, resulting in integer overflow. However, when $\lambda = 0.5, J = 2.5$, the estimator stays
 3329 stable across samples and does not diverge. Still, the estimator does not converge but rather oscillates,
 3330 demonstrating the poor sample complexity of the estimator as well—although it is asymptotically
 3331 unbiased, it may take many samples to converge.
 3332
 3333
 3334

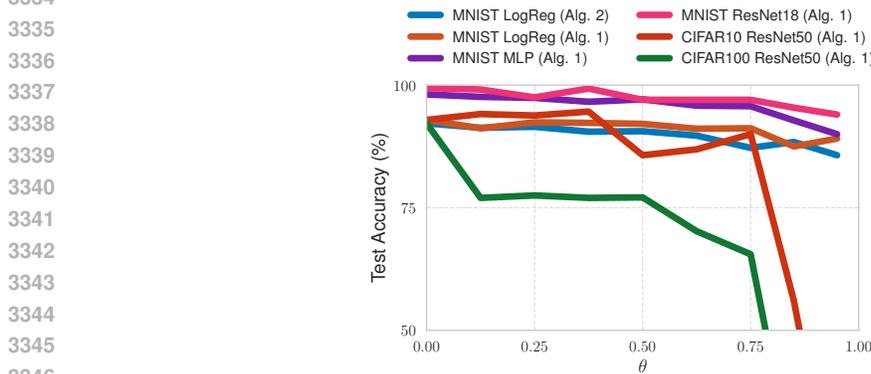
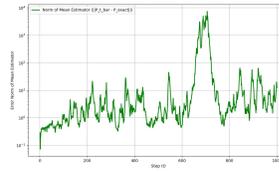
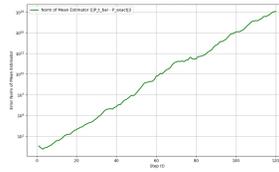
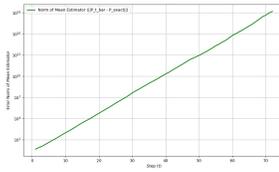


Figure L.12: Test Accuracy vs. θ , MNIST. This has similar behavior to Fig. 3a.

3348 Table L.24: Results for Alg. 3 on MNIST across different models. We compare against Pretraining,
 3349 Alg. 1, and Alg. 2 (where available). We observe that Alg. 3 achieves competitive uniformity (Conf.
 3350 Dist.) on MLP, though with a trade-off in retain accuracy.
 3351

Dataset	Model	Method	Retain Acc.	Test Acc.	Conf. Dist. (Lower Better)
MNIST	LogReg	Pretrain	92.1%	91.8%	0.807
		Alg. 1	87.8%	87.7%	0.180
		Alg. 2	87.1%	87.2%	0.280
		Alg. 3	77.9%	77.3%	0.267
	SVM	Pretrain	92.3%	92.1%	0.543
		Alg. 1	90.3%	90.1%	0.241
		Alg. 3	83.9%	83.9%	0.215
	MLP	Pretrain	98.3%	97.0%	0.886
		Alg. 1	97.5%	95.7%	0.037
Alg. 3		88.8%	87.5%	0.057	



3365
 3366
 3367
 3368
 3369
 3370
 3371 (a) $\lambda = 1.0, J = 2.5$. Here, J is off by 0.5, and the estimator diverges. (b) $\lambda = 0.5, J = 2.0$. Here, λ is off
 3372 by 0.5, and the estimator diverges. (c) $\lambda = 0.5, J = 2.5$. Here, the
 3373 estimator stays roughly stable.

3374 Figure L.13: We observe that small changes in hyperparameters result in divergence of the inverse
 3375 Hessian estimator in Alg. 3.
 3376

3377 M BROADER IMPACTS

3378
 3379 Potential positive impacts are motivated by the threat model as discussed in Sec. 1 and Appx. A ; per
 3380 our example provided in the introduction, violations of test-time privacy constitute a real threat for
 3381 ML safety. Hence, providing the defense that we do constitutes a positive societal impact. However,
 3382 we acknowledge the potential danger in providing a new threat model—it is possible that potential
 3383 adversaries had not thought of this before. Still, we provide a way to address this threat.
 3384
 3385
 3386
 3387
 3388
 3389
 3390
 3391
 3392
 3393
 3394
 3395
 3396
 3397
 3398
 3399
 3400
 3401

3402 N SYMBOL TABLE

3403
3404
3405
3406
3407
3408
3409
3410
3411
3412
3413
3414
3415
3416
3417
3418
3419
3420
3421
3422
3423
3424
3425
3426
3427
3428
3429
3430
3431
3432
3433
3434
3435
3436
3437
3438
3439
3440
3441
3442
3443
3444
3445
3446
3447
3448
3449
3450
3451
3452
3453
3454
3455

Symbols

f	A pretrained classifier.
f_u	A pretrained classifier after unlearning has been conducted over \mathbf{x}_p .
\mathbf{x}_p	A data instance corresponding to person p .
\mathcal{X}	A sample space, subset of \mathbb{R}^d .
\mathcal{Y}	A label space, subset of \mathbb{R}^o .
\mathcal{Z}	The Cartesian product of a sample space and a label space. This is the space where a dataset is drawn from.
\mathcal{D}	A dataset, subset of \mathcal{Z}^n , which is the n-fold Cartesian product of \mathcal{Z} . This represents a set of n data instances.
\mathcal{D}_f	A forget set, a subset of a dataset \mathcal{D} .
\mathcal{D}_r	A retain set, the complement of the forget set in \mathcal{D} .
\mathcal{W}	A space of parameters, subset of \mathbb{R}^p .
A	A function that maps datasets to parameters; this represents the learning algorithm that a ML model provider uses throughout our paper.
$\mathcal{H}_{\mathcal{W}}$	A set of functions which map samples in \mathcal{X} to the probability simplex $\Delta_{ \mathcal{Y} }$, parameterized by a $\mathbf{w} \in \mathcal{W}$.
$\ \mathbf{w}\ _2$	The ℓ_2 norm of a vector \mathbf{w} .
$\ A\ _2$	The 2-operator norm of a matrix A .
λ	An ℓ_2 regularization coefficient used in Alg. 1 for regularization and the certified algorithms e.g. Alg. 2 for local convex approximation.
$\lambda_{\min}(A)$	The minimum eigenvalue of a matrix A .
\mathcal{K}	A uniform learner, which maps samples to parameters which, when one parametrizes a function by any such parameter, a uniform distribution over all possible labels, $U[0, \mathcal{Y}]$ is outputted.
$f_{\mathbf{w}}$	A classifier parameterized by a parameter $\mathbf{w} \in \mathcal{W}$.
\mathcal{L}_A	A loss function to yield accurate predictions e.g. the cross entropy loss between model predictions and labels.
$\mathcal{L}_{\mathcal{K}}$	A loss function to yield accurate uniformity e.g. the Kullback-Liebler divergence between softmax outputs and uniform distribution.
θ	A trade off parameter in $(0, 1)$ between utility and uniformity.

Symbols

\mathcal{M}_θ	A map between datasets and parameters that is the minimizer of a Pareto objective between \mathcal{L}_A and \mathcal{L}_K , where θ spans the (convex) Pareto frontier.
J	The bound on the norm of the Hessian at w^* .
$\mathcal{D}_f^{(i)}$	The i th instance of the forget set.
$\mathcal{D}_r^{(j)}$	The j th instance of the retain set.
$\mathcal{D}_r^{(j,\mathcal{X})}$	The feature of the j th instance of the retain set.
$\mathcal{D}_r^{(j,\mathcal{Y})}$	The label of the j th instance of the retain set.
$\approx_{\varepsilon,\delta,\mathcal{T}}$	Used to denote when two algorithms are (ε, δ) indistinguishable across all subsets $\mathcal{T} \subset \mathcal{W}$, i.e. $\mathcal{M}(\mathcal{D}) \approx_{\varepsilon,\delta,\mathcal{T}} \mathcal{M}'(\mathcal{D}')$ means that $\Pr[\mathcal{M}(\mathcal{D}) \in \mathcal{T}] \leq e^\varepsilon \Pr[\mathcal{M}'(\mathcal{D}') \in \mathcal{T}] + \delta$ and $\Pr[\mathcal{M}'(\mathcal{D}') \in \mathcal{T}] \leq e^\varepsilon \Pr[\mathcal{M}(\mathcal{D}) \in \mathcal{T}] + \delta$.
C	A bound on the model weights.
P_K	The Lipschitz constant for the gradients of ℓ_K , the component loss functions of \mathcal{L}_K , from Asm. 5.3.
P_A	The Lipschitz constant for the gradients of ℓ_A , the component loss functions of \mathcal{L}_A , from Asm. 5.3.
F_K	The Lipschitz constant for the Hessians of ℓ_K , the component loss functions of \mathcal{L}_K , from Asm. 5.4.
F_A	The Lipschitz constant for the Hessians of ℓ_A , the component loss functions of \mathcal{L} , from Asm. 5.4.
P	A convex combination of P_K and P_A with respect to θ .
F	A convex combination of F_K and F_A with respect to θ .
$\mathcal{N}(0, \sigma^2 \mathbf{I})$	The standard normal distribution with an isotropic covariance matrix.
$\nabla_{w^*,\mathcal{K},\mathcal{A}}$	The gradient of the Pareto objective evaluated at w^* , used in the main paper with regularization.
$\mathbf{H}_{w^*,\mathcal{K},\mathcal{A}}$	The Hessian of the Pareto objective evaluated at w^* , used in the main paper without regularization.