
A New Linear Scaling Rule for Differentially Private Hyperparameter Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A major direction in differentially private (DP) machine learning is DP fine-tuning:
2 pretraining a model on a source of public data and transferring the extracted
3 features to downstream tasks. This is an important setting because many industry
4 deployments fine-tune publicly available feature extractors on proprietary data
5 for downstream tasks. In this paper we propose a new linear scaling rule, a
6 hyperparameter optimization algorithm that privately selects hyperparameters to
7 optimize the privacy-utility tradeoff. A key insight into the design of our method
8 is that our new linear scaling rule jointly increases the step size and number of
9 steps as ϵ increases. Our work is the first to obtain state-of-the-art performance
10 on a suite of 20 benchmark tasks across computer vision and natural language
11 processing for a wide range of $\epsilon \in [0.01, 8.0]$ while accounting for the privacy cost
12 of hyperparameter tuning.

13 1 Introduction

14 Industry deployments make use of pretrained
15 models [79] by fine-tuning on task-specific
16 datasets [35; 6; 69] and serving consumer ap-
17 plications that span the range of modalities from
18 portraiture [65] to chatbots [44]. A crucial com-
19 ponent of interfacing machine learning models
20 closely with user data is ensuring that the pro-
21 cess remains *private* [74], and Differential Pri-
22 vacy (DP) is the gold standard for quantifying
23 privacy risks and providing provable guarantees
24 against attacks [20]. DP implies that the output
25 of an algorithm e.g., the final weights trained by
26 stochastic gradient descent (SGD) do not change
27 much if a single datapoint in the dataset changes.

28 **Definition 1.1** (Differential Privacy). A random-
29 ized mechanism \mathcal{M} with domain \mathcal{D} and range \mathcal{R}
30 preserves (ϵ, δ) -differential privacy iff for any
31 two neighboring datasets $D, D' \in \mathcal{D}$ and for
32 any subset $S \subseteq \mathcal{R}$ we have $\Pr[\mathcal{M}(D) \in S] \leq$
33 $e^\epsilon \Pr[\mathcal{M}(D') \in S] + \delta$

34 where D and D' are neighboring datasets if they differ in a single entry, ϵ is the privacy budget and δ
35 is the failure probability.

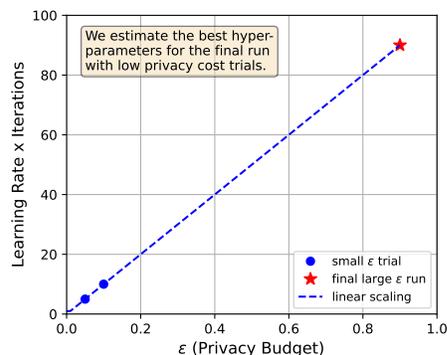


Figure 1: Our new linear scaling rule first does a small number of trials with a very small privacy budget, then does a small number of trials with a slightly larger privacy budget, and finally does linear interpolation through the optimal hyperparameters from these low-cost runs up to the final privacy cost

36 Differentially Private Stochastic Gradient Descent (DP-SGD) [72; 1] is the standard privacy-
 37 preserving training algorithm for training neural networks on private data, with an update rule
 38 given by $w^{(t+1)} = w^{(t)} - \frac{\eta_t}{|B_t|} (\sum_{i \in B_t} \frac{1}{C} \mathbf{clip}_C(\nabla \ell(x_i, w^{(t)})) + \sigma \xi)$ where the changes to SGD are
 39 the per-sample gradient clipping $\mathbf{clip}_C(\nabla \ell(x_i, w^{(t)})) = \frac{C \times \nabla \ell(x_i, w^{(t)})}{\max(C, \|\nabla \ell(x_i, w^{(t)})\|_2)}$, and addition of noise
 40 sampled from a d -dimensional Gaussian distribution $\xi \sim \mathcal{N}(0, 1)$ with standard deviation σ . These
 41 steps alter the bias-variance tradeoff of SGD and degrade utility, creating a challenging privacy-utility
 42 tradeoff. Recent work has made significant progress in closing the gap in performance between
 43 private and non-private fine-tuning of transformer-scale models [46; 52; 7; 51], but a key problem
 44 presents a concrete obstacle to implementing DP algorithms to power real-world consumer-facing
 45 machine learning applications.

46 The privacy analysis of current approaches for private
 47 training does not account for the cost of hyperparameter
 48 tuning, and DP-SGD additionally increases the hyperpa-
 49 rameter tuning burden compared to vanilla SGD. These
 50 hyperparameters include the learning rate schedule, the
 51 clipping bound, the batch size, and the amount of noise to
 52 add at each iteration. Because private training introduces
 53 additional hyperparameters, biases optimization by clip-
 54 ping the gradient, and imposes privacy-utility tradeoffs for
 55 existing hyperparameters, it is challenging to apply hyper-
 56 parameter selection strategies from non-private training,
 57 even on the same dataset. Furthermore prior SOTA work
 58 in private training does not use similar hyperparameters as
 59 non-private training so hyperparameter search algorithms
 60 cannot be leveraged from the broader literature. More
 61 specifically, conventional non-private training uses SGD
 62 with momentum [61] or AdamW [36] to train for hundreds
 63 of epochs. However, training for additional iterations in
 64 DP-SGD requires adding additional noise [27], and taking
 65 large step sizes (such as with momentum) with low signal-
 66 to-noise ratio (SNR) can destabilize training [3]. Prior
 67 work aims to minimize the amount of noise that is added during training by utilizing early stopping,
 68 training for as little as a single iteration [51]. Prior work has either fixed these hyperparameters
 69 without explanation [7] or performed an extensive search to find the best values [15], but the hundreds
 70 of trials of hyperparameter tuning [51] go unaccounted for in the privacy analysis.

71 We propose a new linear scaling rule (Alg. 1, Fig. 1) that automatically selects hyperparameters
 72 to optimize the privacy-utility tradeoff of private fine-tuning. In particular, as our privacy budget
 73 increases from $\epsilon = 0 \rightarrow \infty$, we increase the step size and number of steps. Our method accounts for
 74 the privacy cost of hyperparameter selection by allotting a small portion of the budget to find the best
 75 hyperparameters at $\epsilon \ll 1$ and scaling these up to $\epsilon = 1$. We summarize our contributions:

- 76 • We demonstrate that our new linear scaling rule reduces the computation and privacy cost of
 77 hyperparameter optimization by an order of magnitude without sacrificing performance
- 78 • Linear scaling can obtain new SOTAs for both full fine-tuning and linear probing of both
 79 convolutional and transformer architectures across 20 vision and language tasks
- 80 • We compare four model architectures for a set of five vision benchmarks and find that the private-
 81 non private utility gap decreases as models improve, with the best model across all five tasks
 82 obtaining lossless performance of 99% accuracy for $\epsilon = 1$ on CIFAR10
- 83 • We find that linear scaling is robust to domain shifts between private training and test data
- 84 • We find that models trained with our method can provide good performance even when there is a
 85 large shift between public and private data
- 86 • We validate that models trained with our method can perform well for zero-shot classification
- 87 • We [provide our code](#) as a part of our empirical evaluation.

Figure 2: We compare the best private and best non-private test accuracy performances of our method to prior work using models pretrained on ImageNet-21k and fine-tuned on CIFAR10 and CIFAR100 datasets. Our results at $\epsilon = 1$ include the cost of hyperparameter tuning via applying the linear scaling rule at $\epsilon \in [0.01, 0.1]$.

Dataset	Approach	$\epsilon = 1$	$\epsilon = \infty$
CIFAR10	Ours	99.00	99.00
	[51]	96.30	96.60
	[7]	96.70	97.40
	[9]	95.00	96.40
	[15]	94.80	96.60
CIFAR100	Ours	89.62	91.57
	[51]	82.70	85.29
	[7]	83.00	88.40
	[9]	73.70	82.10
	[15]	67.40	81.80

Algorithm 1 DP-SGD with Linear Scaling

Inputs: Private dataset \mathcal{D} , open source feature extractor F , number of classes C , privacy budget ε , momentum $\rho = 0.9$, first search privacy budget ε_0 , second search privacy budget ε_1
Perform first hyperparameter search to obtain the best possible value of r_0 within the first privacy budget ε_0
Perform second hyperparameter search initialized at $r_1^* = \frac{\varepsilon_1}{\varepsilon_0} \cdot r_0$ to obtain the best possible value of r_1 within the second overall privacy budget ε_1
Perform linear interpolation to estimate the slope α and bias b of the line $r = \alpha\varepsilon + b$ given $(r_0, \varepsilon_0), (r_1, \varepsilon_1)$
Set $r^* = \alpha\varepsilon_f + b$ given the estimated linear interpolation
Extract features from \mathcal{D} using F : $\mathcal{X} = F(\mathcal{D})$
Zero-initialize classifier $w \leftarrow 0_{C \times d}$
Decompose the total step size r given by linear scaling into $r = \eta \times T$
Use privacy loss variable accounting to calibrate noise parameter σ given ε
for $i = 1, 2, \dots, T$ **do**
 Compute full-batch gradient according to Eq. 1 $\nabla^{(i)} = \frac{1}{|\mathcal{D}|} (\sum_{i \in \mathcal{D}} \text{clip}_1(\nabla \ell(x_i, w^{(i)})) + \sigma \xi)$
 Take a step with momentum: $v^{(i)} \leftarrow \rho \cdot v^{(i-1)} + \nabla^{(i)}, w^{(i)} \leftarrow w^{(i-1)} - \eta v^{(i)}$
end for
Output: $(\varepsilon_f + \varepsilon_0 + \varepsilon_1)$ -Private linear model w

88 2 A New Linear Scaling Rule

89 In this section we detail how our method chooses each hyperparameter in DP-SGD, prove the privacy
90 guarantee of the overall hyperparameter selection process, and provide a theoretical analysis.

91 **A new linear scaling rule** The well-known linear scaling rule [29] proposes increasing the learning
92 rate with the batch size. We propose a new linear scaling rule that details how to select all hyperpa-
93 rameters in DP-SGD. Our method first fixes full-batch, unit clipping norm, zero initialization and
94 use SGD with momentum, and then jointly scales the learning rate and number of steps with ε . We
95 provide extensive ablations of each design choice in our hyperparameter optimization algorithm
96 in Appendix A.2. Prior work has exclusively taken small step sizes [51; 52; 7; 15; 9] on the order
97 of $\{10^{-5}, 10^{-3}\}$ and works that train transformers have also trained for a small number of epochs
98 $\{1, 3\}$ [51; 7]. While this works well to recover the bulk of the non-private performance when ε is
99 very small, it is natural to expect that as $\varepsilon \rightarrow \infty$ we should increase the parameters of training to more
100 closely resemble that of non-private training. In line with this insight, we propose a linear scaling
101 rule: jointly increase the step size and number of steps linearly with ε . We make use of this simple
102 yet powerful heuristic in the hyperparameter selection strategy that we use in all our experiments,
103 outlined in Algorithm 1. Given a total privacy budget ε , we use an initial portion of this budget to do
104 binary search (random search and grid search are also valid) on the meta-hyperparameter $r = \eta \times T$
105 for a small value of ε , and use this to estimate the best value of r for the desired overall privacy
106 budget. We provide a privacy guarantee in 2. We note that linear scaling does not hold up forever: we
107 are primarily interested with analyzing $\varepsilon \leq 1$, and show that in this range it holds (Fig. 3).

108 **Linear Scaling is intuitive.** Applying the linear scaling rule improves the cosine similarity between
109 noisy weight updates and the optimal solution without degrading accuracy. First note that the
110 classification accuracy of a linear model is scale-invariant; the optimal solution of Gradient Descent
111 with total step size r is $w' = w^* / \|w^*\| \times r$: the projection of w^* onto B_r , the ball of radius r ,
112 and for linear models, the performance (top-1 accuracy) of w' is the same as the performance of
113 w^* : $\text{Pred}(w'(x)) = \text{Pred}(w^*(x)) \forall x \in \mathcal{D}$. An important factor in the success of optimization is
114 the angle between the gradient update ∇_i and w' : if all our updates point in the same direction, we
115 can expect fast convergence. Let similarity(i) = $\frac{\nabla_i \cdot w'}{\|\nabla_i\| \cdot \|w'\|}$. Suppose that $\|w_i\| = \|w'\| \ll 1$, then
116 adding Gaussian noise $\sigma \xi$ where $\xi \sim \mathcal{N}(0, 1)$ to the update will significantly decrease the cosine
117 similarity of the updated model and w' . If we decrease σ , it is easy to see that this mitigates the
118 impact on the trajectory. However, we can equivalently keep σ constant and increase the scale of
119 the parameters, and also decrease the impact of noise on the trajectory: similarity($w_i + \sigma \xi, w'$) <
120 similarity($\alpha \cdot w_i + \sigma \xi, \alpha \cdot w'$), $\forall \alpha > 1$. Note that by increasing r we scale the optimal solution

121 while keeping its performance identical, and thus optimize the cosine similarity of the noisy update.
 122 Increasing the number of iterations and the learning rate linearly increases r but does not linearly
 123 increase σ due to the composition of Gaussian differential privacy [27], therefore the impact on the
 124 optimization trajectory is minimized.

125 **Theory** We introduce two theoretical results. We first analyze the privacy cost including hyperpa-
 126 rameter tuning of DP-RAFT under Gaussian DP (GDP). In Thm. 2.3 we analyze the performance
 127 gap between hyperparameters for noisy gradient descent in terms of an upper bound in expectation
 128 on the distance between private and non-private iterates, and find that applying the linear scaling rule
 129 improves the upper bound on this distance. Proofs of all results are in Appendix A.5.

130 **Proposition 2.1.** *Algorithm 1 is (\sqrt{T}/σ) -GDP. Moreover, repeating Algorithm 1 for n times for
 131 hyper parameter search would be $(\sqrt{T} \cdot n/\sigma)$ -GDP.*

132 **Corollary 2.2.** *Algorithm 1 is $(\epsilon, \Phi(-\epsilon \cdot \sigma/\sqrt{T} + \sqrt{T}/2\sigma)) - e^\epsilon \cdot \Phi(-\epsilon \cdot \sigma/\sqrt{T} - \sqrt{T}/2\sigma)$ -DP. Also,
 133 for n -fold repetition, the algorithm is $(\epsilon, \Phi(-\epsilon \cdot \sigma/\sqrt{n \cdot T} + \sqrt{n \cdot T}/2\sigma)) - e^\epsilon \cdot \Phi(-\epsilon \cdot \sigma/\sqrt{n \cdot T} -$
 134 $\sqrt{n \cdot T}/2\sigma)$ -DP*

135 **Theorem 2.3.** *Let f be gradient descent that minimizes a α -strongly convex and β -smooth function
 136 ℓ with constant learning rate $\eta \in (0, \frac{2}{\beta})$ over T iterations. Then we can bound the "noisy radius"
 137 distance between the noisy iterate w^T and the benign iterate w_b^T at iteration T in expectation:
 138 $\mathbb{E}[|w^T - w_b^T|] \leq \rho\eta \times (\sum_i^{T-1} \max(|1 - \eta\alpha|, |1 - \eta\beta|)^i)$.*

139 Thm. 2.3 indicates that the distance between the noisy and non-noisy weights grows in a very
 140 controlled manner; at each iteration the divergence from the previous iteration is decreased by a factor
 141 strictly less than 1, and then we add some noise. The main idea of the proof is similar to the main
 142 result in Fang et al. [23] but is simpler because we only prove the result for linear models.

143 We apply this theorem to logistic regression (fine-tuning a linear model on extracted features). In this
 144 setting our theorem provides an upper bound on the radius of the range of solutions that DP-SGD
 145 produces. For linear models, this radius converts directly into an upper bound on the generalization
 146 error. If we use the linear scaling rule to scale $r = \eta \times T$ with ϵ , we expect that η remains
 147 appropriately bounded and T does not grow so large that the resulting noise creates significant model
 148 drift. Therefore, we find that increasing the quantity $r = \eta \times T$ improves this upper bound.

149 While our theorem only holds for linear models, we will show that it holds empirically for the deep
 150 GPT2 and RoBERTa models, in line with Li et al. [47] who find that even the updates of a large
 151 model lie in a low-dimensional space during fine-tuning.

152 3 Evaluation

153 We provide results on a range of image classification, distribution shift, and natural language process-
 154 ing tasks. Full results for all datasets and models can be found in Appendix A, including ablations on
 155 all steps of our method(A.2) and key hyperparameters(A.4).

156 **Datasets.** We evaluate the performance of our method on 20 benchmark tasks spanning the data
 157 modalities of CV and NLP. Image classification: ImageNet [16], CIFAR10, CIFAR100 [40], Fashion-
 158 MNIST [80], STL10 [11], EMNIST [12]. Because these image classification datasets are generally
 159 considered in-distribution of the pretraining data, we also provide results on a number of distribu-
 160 tion shift datasets from the WILDS suite [38] that have been used to evaluate various fine-tuning
 161 techniques. CIFAR10 \rightarrow STL, CIFAR10p1, CIFAR10C, CIFAR100 \rightarrow CIFAR100C [31], Water-
 162 birds [67], FMoW [10], and Camelyon17 [8]. These datasets are considered benchmark tasks for
 163 distribution shifts [42; 43; 53] and include data that is not in-distribution of the training data, making
 164 for a more realistic evaluation of the capabilities of our method to solve challenging tasks. We are
 165 the first to show that DP-SGD is capable of learning to handle distribution shifts without using any
 166 techniques from the distributionally robust optimization (DRO) literature [64]. For NLP tasks we
 167 consider text classification tasks from the GLUE benchmark [76]: SST-2, QNLI, QQP, MNLI(m/mm)
 168 and for next word generation we use PersonaChat [84], WikiText-2 [54], and Enron Emails [37].

169 **3.1 Linear Scaling finds near-optimal hyperparameters with low privacy cost**

170 We first provide a concrete example of the hyperparameter search with ε_0 on CIFAR10. Note that
 171 regardless of what strategy we use for hyperparameter search here, our total privacy cost as given
 172 by Proposition 2 must be strictly less than ε_0 . Binary search, random search, Bayesian optimization
 173 and grid search are all methods that we can use for the initial hyperparameter search. For this
 174 example, for the sake of simplicity we will use random search with 3 trials, with $\varepsilon_0 = 0.01 \cdot \sqrt{3}$, $\varepsilon_1 =$
 175 $0.05 \cdot \sqrt{3}$, $\varepsilon_f = 0.9$, $\varepsilon_0 + \varepsilon_1 + \varepsilon_f = 1.0$. For $\varepsilon_0 = 0.01$, we randomly sample r uniformly in the
 176 range $[1, 100] = 2, 20, 100$ and then randomly decompose this into (approximate) (η, T) pairs of $[0.2,$
 177 $10]$, $[0.5, 40]$, $[1, 100]$. These in turn evaluate to accuracies of $[91.79, 73.68, 67.21]$, so the best value
 178 of r at $\varepsilon_0 = 0.01$ is 2. We do a similar process at $\varepsilon_1 = 0.05$ and get a best r -value of 5. We do linear
 179 interpolation and obtain the line of best fit as $r = 75 \cdot \varepsilon + 1.25$. Approximating this to $r = 75$, we
 180 apply the linear scaling rule $r = \eta \times T$ and randomly decomposing this value of r into an (η, T) pair
 181 of $[0.75, 100]$, we produce a final accuracy of 99.00 at $\varepsilon_f = 0.9$.

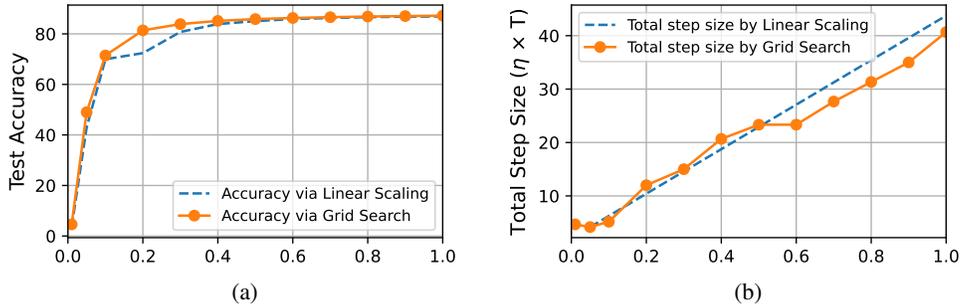


Figure 3: Training the beit architecture on CIFAR100, the linear scaling rule produces values for $r = \eta \times T$ close to that of grid search, and the performance drop is only apparent at $\varepsilon > 0.2$ because the cost of tuning is $\varepsilon = 0.1$, and vanishingly small for larger ε .

182 **Linear Scaling outperforms prior hyperparameter search techniques.** We validate the effective-
 183 ness of linear scaling against the grid search baseline. In Fig. 3 (right) we compare Alg. 1 to
 184 grid search. To avoid scale mismatch on the x-axis we do not account for the privacy cost of grid
 185 search, that does $n = 100$ trials (on the same scale as prior work [51]). It is trivial that linear scaling
 186 outperforms a naive grid search, but we also compare the effectiveness of linear scaling against
 187 the hyperparameter selection strategies used in prior work [51]. We apply linear scaling to the ViT
 188 model used in [51] on CIFAR100. Although [51] do not directly state the hyperparameters for
 189 their best results, they specify that they use 200 hyperparameter trials with Bayesian optimization.
 190 While they obtain RDP guarantees, these guarantees do not include the privacy cost of non-privately
 191 tuning hyperparameters. We apply the linear scaling rule to extrapolate a value of r from $\varepsilon = 0.1$
 192 to $\varepsilon = 1$, obtaining $r = 20 = \eta(0.2) \times T(100)$. We recover performance of 82.7% for $\varepsilon = 1$, a
 193 2% improvement over the best result for DP-Adam in [51] while accounting for the privacy cost
 194 of hyperparameter tuning. They obtain their best result for DP-Adam at $T = 10$, but we cannot
 195 compute the corresponding value of r because they do not provide η . However, because they use a
 196 clipping norm of 0.005 we can reasonably infer that their value of r is $\approx 1000\times$ smaller than ours.
 197 This is farther from the optimal non-private training, as evidenced by the performance gap.

198 Figure 4: Linear Scaling on ImageNet is competitive with prior SOTA [52] (Jan. 2023) and
 199 current SOTA [51](within last month).
 200

ε	[52]	[51]	Ours	$r = \eta \times T$
0.25	75.6	-	79.0	250
0.50	79.4	86.1	81.6	750
1.00	81.1	86.8	83.2	1100
2.00	81.5	87.4	84.2	2000
10.0	81.7	-	85.4	2000
∞	86.9	88.9	85.7	2000

Linear Scaling scales to ImageNet In Table 4
 we do a granular comparison between our method
 and [52; 51]. We observe that our method is compet-
 itive with [51] even when accounting for the privacy
 cost of hyperparameter search, and that the linear
 scaling rule holds up at the scale of ImageNet for
 very large values of $r = \eta \times T$. The non-private
 accuracy of their closed-source model is 3.2% higher
 than our open-source model, and so the private
 accuracy at $\varepsilon = 2$ is also 3.2% higher.

208 However, ultimately our method and the method of Mehta et al. [51] are complementary, because
 209 their method introduces new hyperparameters that we intuit our linear scaling rule can optimize. We
 210 attempted to validate this intuition empirically but were unable to reproduce the results of Mehta et al.
 211 [51] because they and Mehta et al. [52] pretrain on the closed-source JFT dataset with billions of
 212 images. We note that all numbers we report for models pretrained on ImageNet-21k using first-order
 213 methods surpass those in [51], but for sufficiently small values of ϵ on harder datasets the second-
 214 order methods they propose provide better performance. We note that the method in Mehta et al. [51]
 215 only works for vision tasks, whereas our approach works for both vision and language tasks.

216 **Linear Scaling produces robust results.** In Fig. 3 we report that following Algorithm 1 produces
 217 new state-of-the-art results for all values of ϵ , shown in Table 5. In Appendix A.1 we provide detailed
 218 computations of the linear interpolation for multiple datasets and in Appendix A.4 we provide full
 219 results across the entire hyperparameter search space. Our results validate that this rule is robust: we
 220 can move from one set of hyperparameters to another similarly performing set of hyperparameters by
 221 increasing the number of iterations T by a constant factor and decreasing the learning rate η by the
 222 same factor (or vice versa). We find that any inaccuracy incurred by estimating the best value of r
 223 with the linear scaling rule will not reduce accuracy by much compared to doing grid search for the
 224 optimal value of r , but does reduce the privacy cost of hyperparameter tuning immensely.

225 3.2 Linear Scaling enables empirical analysis

226 Many interesting questions in DP fine-tuning remain unanswered because of the immense compu-
 227 tational overhead of evaluating hundreds of hyperparameter trials for each privacy budget, model
 228 architecture and dataset [51]. We now employ the linear scaling rule to efficiently answer key
 229 questions in DP fine-tuning for vision tasks.

230 **Impact of model architectures on differential**
 231 **privacy** Many pretrained model architectures
 232 are available [79] but prior work has generally
 233 engaged with a single architecture, e.g. beit [7]
 234 or ViT [52]. We leverage our method to answer
 235 three questions:

- 236 • What model architectures can provide good
 237 DP classifiers?
- 238 • Is the best model task-specific, e.g., is an
 239 architecture search required?
- 240 • Does the private-non private utility gap de-
 241 pend on the model architecture?

242 We report our findings in Tab. 5. We evaluate
 243 multiple transformer architectures in ViT [19],
 244 beitv1 [4] and beitv2 [58], as well as the purely
 245 convolutional architecture Convnext [48]. We
 246 find that all architectures can serve as good back-
 247 bones for high-accuracy DP classification. This
 248 is somewhat surprising because the different in-
 249 ductive biases of transformers and purely convo-
 250 lutional architectures tend to produce differently

251 structured features, but we reason that the noise added by DP will ‘smooth out’ these decision
 252 boundaries regardless of architecture. We note that one architecture, beitv2, performs the best on all
 253 benchmarks and also has the highest non-private ImageNet accuracy [78]. We therefore recommend
 254 that practitioners do not worry about architecture search when fine-tuning as this can incur further
 255 privacy costs, and instead pick the best model available. We are encouraged to report that the
 256 private-non private utility gap diminishes with model accuracy, enabling us to report for the first time
 257 *lossless privacy* of 99.0% on CIFAR10 at $\epsilon = 1$. We expect that as pretrained models become even
 258 better, future works may even be able to attain lossless privacy on CIFAR100, that we note remains
 259 somewhat challenging for private fine-tuning. We harness these insights for our next analyses.

Figure 5: We compare the best private and best non-private performances of all models on all datasets. We use the linear scaling rule to scale hyperparameters from $\epsilon = 0.1$ to $\epsilon = 1$, so our privacy analysis includes the cost of hyperparameter tuning.

Model	Dataset	$\epsilon = 1$	$\epsilon = \infty$	Gap
beitv2	CIFAR10	99.00	99.00	0.00
	CIFAR100	89.62	91.57	1.95
	FMNIST	91.02	91.53	0.51
	STL10	99.69	99.81	0.12
	EMNIST	81.77	82.00	0.23
convnext	CIFAR10	96.75	97.22	0.47
	CIFAR100	83.47	86.59	3.12
	FMNIST	90.23	91.13	0.9
	STL10	99.61	99.71	0.10
	EMNIST	78.38	79.05	0.67
beit	CIFAR10	98.19	98.51	0.32
	CIFAR100	87.1	90.08	2.98
	FMNIST	90.55	91.6	1.05
	STL10	99.62	99.78	0.16
	EMNIST	81.48	83.25	1.77
vit-L	CIFAR10	98.29	98.44	0.40
	CIFAR100	86.18	89.72	3.54
	FMNIST	90.58	91.37	0.79
	STL10	99.62	99.76	0.14

260 **Linear Scaling is robust to distri-**
 261 **bution shifts.** Benchmarking per-
 262 formance on datasets with distri-
 263 bution shifts is increasingly impor-
 264 tant because real-world problems al-
 265 most always contain distribution shift
 266 between model training and infer-
 267 ence [64]. Prior work in distribution-
 268 ally robust optimization (DRO) has ad-
 269 dressed this problem by using knowl-
 270 edge of the relative imbalances be-
 271 tween groups, but recent work with
 272 vision transformers has shown that
 273 linear probing can perform well on
 274 datasets with distribution shifts [53;
 275 41; 43]. However there is no work
 276 that evaluates the robustness of private models to distribution shifts. We leverage our method to
 277 answer three questions:

- 278 • Can DP help when there is a domain shift from private fine-tuning to test?
- 279 • Can DP help when there is a domain shift from public data to private fine-tuning?
- 280 • Can DP fine-tuned models perform well in the zero-shot setting?

281 In Table 6 we compare the performance of our method across 8 benchmarks and find that the answer
 282 to all three of these questions is *yes*.

283 The Waterbirds dataset is a well-known benchmark for evaluating the robustness of models to spurious
 284 correlations. There is a domain shift between the private training data and the private test data created
 285 by class imbalance. We are surprised to find that in the absence of any other regularization methods,
 286 DP fine-tuning actually *improves* performance on the OOD split. We hypothesize that the lackluster
 287 OOD non-private performance is caused by the model overfitting to the spurious correlation in the
 288 training data, and that the inherent regularization of DP prevents the model from memorizing this
 289 spurious correlation. By comparing our results to Mehta et al. [53] we determine that this robustness
 290 is unique to DP rather than an artifact of the pretrained model. Although DP does significantly
 291 degrade the ID performance, in situations where minimizing OOD error is more important, we believe
 292 that DP by itself can mitigate the domain shift from private fine-tuning to test.

293 Because our central assumption in DP fine-tuning is that there is no privacy leakage from the
 294 pretraining data to the private training data, it is important to understand how DP fine-tuning
 295 performs when there is a distribution shift between public data and private data. fMoW [10] and
 296 Camelyon17 [8] are two datasets that represent a significant distribution from the pretraining data
 297 (ImageNet). We observe a similar relationship between ID and OOD degradation as above, where the
 298 OOD degradation is somewhat mitigated by DP. If we compare our results on Camelyon to the best
 299 results in Ghalebikesabi et al. [25] we find that we can improve their best performance from 91.1% at
 300 $\epsilon = 10$ to 93.91% at $\epsilon = 1$. Although performance on fMoW remains quite poor, we note that it is
 301 not significantly worse than in the non-private setting. We believe that DP fine-tuning from pretrained
 302 models remains a viable strategy even when the publicly available pretraining data has a very large
 303 distribution shift from the private target data.

304 We finally consider the zero-shot setting, where we fine-tune a model on CIFAR and then transfer it
 305 without updating any parameters to private test datasets that once again represent a distribution shift
 306 from CIFAR. We report the performance in the OOD column. For the more minute distribution shifts
 307 of STL and CIFAR10p1 we find that the fine-tuned classifier can achieve remarkable performance
 308 without ever updating parameters on these datasets; that is, we just remap the labels as per [42].
 309 CIFAR10C and CIFAR100C represent larger distribution shifts and are used to benchmark the
 310 robustness of models to commonly reported image corruptions [31]. Our OOD performance on
 311 these larger distribution shifts is much worse, particularly for CIFAR100 where there is a $> 20\%$
 312 degradation. Although this is lower than the top result on the RobustBench leaderboard [13] obtains
 313 85% accuracy, we note that once again *we used no additional methods beyond DP to ensure robustness*
 314 *but managed to achieve reasonable performance to distribution shifts in zero-shot classification.*

Figure 6: In-distribution (ID) and out-of-distribution (OOD) performance on benchmark distribution shift datasets. Prior work is non-private (citations are in Appendix A.1). We use the linear scaling rule to scale hyperparameters from $\epsilon = 0.1$ to $\epsilon = 1$, so our privacy analysis includes the cost of hyperparameter tuning.

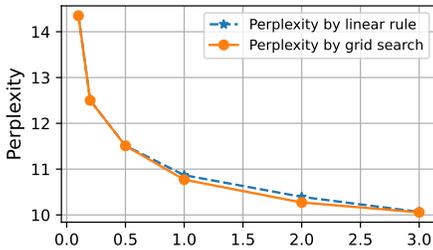
Dataset	$\epsilon = 1.0$ ID(OOD)	Prior ($\epsilon = \infty$)
Waterbirds	92.31 (91.59)	98.3(80.4)
fMoW	45.44 (35.31)	49.1 (36.6)
Camelyon	93.91 (93.55)	99.5 (96.5)
C10 \rightarrow STL	99.0 (98.82)	97.5 (90.7)
C10 \rightarrow C10p1	99.0 (97.85)	97.5 (93.5)
C10 \rightarrow C10C	99.0 (89.98)	96.56 (92.78)
C100 \rightarrow C100C	89.65 (68.69)	81.16 (72.06)

3.3 Linear Scaling for language modeling

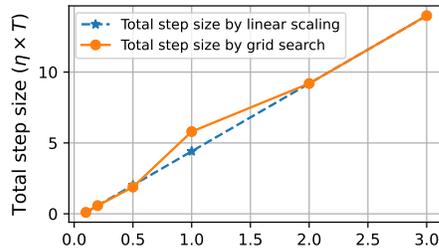
Prior work has generally focused on either CV or NLP because the methods used in DP fine-tuning differ greatly across data modalities [46; 51]; here we show that our method extends to NLP by validating on text classification and language modeling tasks. We also update all parameters when fine-tuning, displaying that our method works for both linear probing and full fine-tuning. We fine-tune GPT-2 [63] with our method for three language modeling tasks that have been benchmarked in prior works [46; 70; 30] on private fine-tuning: Persona-Chat [85], WikiText-2 [54] and Enron Emails [37]. We also fine-tune RoBERTa-base on four tasks in the GLUE benchmark: SST-2, QNLI, QQP and MNLI(m/mm) in Table 7.

While prior works mainly focus on ϵ in $\{3, 8\}$, in this work we are also interested in smaller ϵ s like 0.1. Appendix B.1 includes the details for the experimental set-up.

Linear scaling holds for NLP tasks We analyze the performance gap between estimated total step size and optimal total step size by grid search to understand how well linear scaling performs on language modeling tasks. Fig. 8 plots the optimal perplexity and perplexity by estimated total step size at different values of ϵ on Enron emails. We can see that the linear scaling rule generalizes well for reported values of ϵ and the perplexity by the estimated total step size is close to the optimal perplexity. From Table 7 we can see that linear scaling also holds across a range of tasks in the GLUE benchmark. We also have the result for WikiText-2 in Appendix B.3.



(a) Pareto Frontier for ϵ vs Test Perplexity.



(b) Pareto Frontier for ϵ vs Total Step Size.

Figure 8: The linear scaling rule (accounting for the privacy cost of hyperparameter tuning) is competitive with grid search (non-private, doing N trials each with the given ϵ) on the Enron Emails dataset. Left: y-axis is Perplexity (lower is better).

The linear scaling rule outperforms prior results on differentially private language modeling tasks. We first run a qualitative evaluation on the previous benchmark SOTA [46] on PersonaChat trained with DP-SGD by following the linear scaling rule to increase the number of epochs.

Figure 9: Linear scaling holds when fine-tuning all layers of GPT2 on PersonaChat and outperforms Li et al. [46]

ϵ ($\delta = \frac{1}{2^{\lfloor D_{\text{train}} \rfloor}}$)	1	3	∞
Li et al. [46]	-	24.59	18.52
Our Work	21.25	-	17.69

341

342

343

344

345

346

347

348

349

350

351

352

353

354

We can see in Table 9 that we can push the perplexity under 18 for $\epsilon = 3$ and $\epsilon = 8$; this performance is competitive with the non-private baseline. Furthermore, even when pushing for a stricter privacy guarantee $\epsilon = 0.5$, we can still get perplexity of 21.25, that is better than the result of $\epsilon = 8$ in [46]. We also report the results of ablating these hyper-parameters and varying the number of layers trained in Appendix B.2.

We quantitatively validate the linear scaling rule on WikiText-2 and Enron email dataset and report the result in Table 10 respectively. We select training parameters and the total step size with Alg. 1.

Figure 7: Linear scaling holds for GLUE tasks when training the full RoBERTa-base model

Task	ϵ	Acc	$r = \eta \times T$
SST-2	0.1	90.60	0.975
	0.2	90.83	1.95
	1.0	91.51	9.75
QNLI	0.1	82.54	3.9
	0.2	84.00	4.68
	1.0	86.25	26.52
QQP	0.1	81.07	11.7
	0.2	82.21	17.55
	1.0	84.69	64.35
MNLI(m/mm)	0.1	77.52(78.24)	11.7
	0.2	79.40(79.98)	17.55
	1.0	81.86(82.76)	64.35

355 For WikiText-2, a key observation is
 356 that when we compare our results to
 357 the best prior reported results in [70],
 358 for the same number of passes over
 359 the training data (20), we obtain lower
 360 perplexity for $\epsilon = 0.2$ than they re-
 361 port for $\epsilon = 3$. That is, by just in-
 362 creasing the effective step size from
 363 $\sim 8 \times 10^{-6}$ to $\sim 8 \times 10^{-3}$ we can
 364 strengthen the privacy guarantee with-
 365 out degrading performance.

Figure 10: Finetuning GPT-2 on WikiText-2 ($\delta = 10^{-6}$) and Enron ($\delta = \frac{1}{2|D_{\text{train}}|}$) with DP-SGD. Ppl is perplexity and TSS is Total Step Size. (* means estimated). Previously reported best perplexity of GPT-2 on WikiText-2 at $\epsilon = 3$ is 28.84 in [70].

Dataset	ϵ	0.1	0.2	0.5	1.0	2.0	3.0
WikiText-2	Ppl	-	28.81	28.37	28.15	27.98	27.69
	TSS	-	0.008	0.02	0.04*	0.08*	0.12*
Enron	Ppl	14.35	12.50	11.56	10.91	10.45	10.14
	TSS	0.10	0.58	2.02*	4.41*	9.19*	13.98*

366 4 Related Work and Discussion

367 De et al. [15] and Cattan et al. [9] propose the use of large batch sizes and initializing the weights to
 368 small values near-zero to standardize training. However, they use ResNet architectures rather than
 369 modern vision transformers, and in Appendix A.2 we find that other techniques that they use such as
 370 data augmentation, fine-tuning the embedding layer, and weight averaging do not always improve
 371 performance. [7] do end-to-end training of the same beit architecture we use, but we crucially
 372 observe that updating all parameters incurs the curse of dimensionality and therefore it is better to
 373 only update the last layer. Besides vision tasks, Li et al. [46] and Yu et al. [82] provide methods for
 374 fine-tuning large language models under DP-SGD by proposing new clipping methods to mitigate
 375 the memory burden of per-sample gradient clipping. However, they do not achieve performance
 376 comparable to non-private models when fine-tuning a pretrained model on the PersonaChat dataset.
 377 We adapt their techniques to the hyperparameter settings that we show are optimal for DP fine-tuning,
 378 and produce similar performance to non-private fine-tuning on the PersonaChat dataset. Yu et al.
 379 [83] report compelling results by only updating a sparse subset of the LLMs with LoRA [33]. We
 380 fine-tune GPT2 and RoBERTA; Basu et al. [5] also fine-tune BERT models.

381 Papernot and Steinke [57] propose an RDP hyperparameter optimization algorithm that requires
 382 selecting the number of trials at random with a random variable, and exhibits the greatest savings
 383 when the number of hyperparameter trials is large. By contrast our linear scaling rule needs only a
 384 small fraction of the overall privacy budget for hyperparameter search. Their evaluation only tunes
 385 the learning rate of a 3-layer CNN on MNIST. Our rule accounts for multiple hyperparameters (batch
 386 size, clipping norm, momentum, learning rate, number of iterations) and produces SOTA results.

387 Golatkar et al. [26]; Nasr et al. [55]; Amid et al. [2] treat $< 10\%$ of the private training dataset and
 388 public and use it to improve DP-SGD. Although we do not use any private data during pretraining,
 389 future work can tackle applying linear scaling to this alternate threat model.

390 An open challenge in DP training is how to privately and efficiently do hyperparameter tuning. We
 391 complement the existing body of work by introducing a new linear scaling rule to privately optimize
 392 hyperparameters. Our key insight is that we can interpolate between the early-stopping regime that
 393 is best for small ϵ and the regime of many iterations that is best for $\epsilon \rightarrow \infty$ as ϵ increases. We
 394 provide find that our method attains new state-of-the-art accuracy across 20 tasks, on benchmark
 395 image classification tasks, distribution shift datasets, and natural language modeling tasks.

396 5 Limitations

397 **Assumptions.** The key assumption in DP fine-tuning is that there is no privacy leakage between
 398 public data and private data. We take steps towards qualifying this assumption by evaluating on
 399 datasets with distribution shifts between public and private data. **Scope of Claims.** We evaluate 20
 400 datasets across multiple data modalities with multiple model architectures for two types of fine-tuning
 401 methods, linear probing and end-to-end training of deep ($> 100M$ param) transformers. **Key Factors**
 402 **that Influence the Performance of Our Approach.** The key parameter in the linear scaling rule
 403 is how to allocate privacy budget to the initial hyperparameter search. We find that with privacy
 404 budgets as small as $\epsilon = 0.01$ we can still effectively forecast the linear trend to determine the best
 405 hyperparameters for the main privacy budget we consider $\epsilon = 1$. However, if we need to consider
 406 even smaller privacy budgets, it may be challenging to accurately extrapolate hyperparameters.

407 **References**

- 408 [1] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep
409 learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on*
410 *Computer and Communications Security*. ACM, oct 2016. doi: 10.1145/2976749.2978318.
- 411 [2] E. Amid, A. Ganesh, R. Mathews, S. Ramaswamy, S. Song, T. Steinke, T. Steinke, V. M.
412 Suriyakumar, O. Thakkar, and A. Thakurta. Public data-assisted mirror descent for private
413 model training. In *Proceedings of the 39th International Conference on Machine Learning*,
414 pages 517–535. PMLR, 2022.
- 415 [3] E. Bagdasaryan and V. Shmatikov. Differential privacy has disparate impact on model accuracy,
416 2019. URL <https://arxiv.org/abs/1905.12101>.
- 417 [4] H. Bao, L. Dong, S. Piao, and F. Wei. Beit: Bert pre-training of image transformers, 2021. URL
418 <https://arxiv.org/abs/2106.08254>.
- 419 [5] P. Basu, T. S. Roy, R. Naidu, Z. Muftuoglu, S. Singh, and F. Mireshghallah. Benchmarking
420 differential privacy and federated learning for bert models, 2022.
- 421 [6] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan,
422 P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child,
423 A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray,
424 B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei.
425 Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- 426 [7] Z. Bu, J. Mao, and S. Xu. Scalable and efficient training of large convolutional neural networks
427 with differential privacy. *arXiv preprint arXiv:2205.10683*, 2022.
- 428 [8] P. Bándi, O. Geessink, Q. Manson, M. Van Dijk, M. Balkenhol, M. Hermsen, B. Ehteshami Be-
429 jnordi, B. Lee, K. Paeng, A. Zhong, Q. Li, F. G. Zanjani, S. Zinger, K. Fukuta, D. Komura,
430 V. Ovtcharov, S. Cheng, S. Zeng, J. Thagaard, A. B. Dahl, H. Lin, H. Chen, L. Jacobs-
431 son, M. Hedlund, M. Çetin, E. Halıcı, H. Jackson, R. Chen, F. Both, J. Franke, H. Küsters-
432 Vandeveld, W. Vreuls, P. Bult, B. van Ginneken, J. van der Laak, and G. Litjens. From
433 detection of individual metastases to classification of lymph node status at the patient level: The
434 camelyon17 challenge. *IEEE Transactions on Medical Imaging*, 38(2):550–560, 2019. doi:
435 10.1109/TMI.2018.2867350.
- 436 [9] Y. Cattan, C. A. Choquette-Choo, N. Papernot, and A. Thakurta. Fine-tuning with differential
437 privacy necessitates an additional hyperparameter search, 2022. URL <https://arxiv.org/abs/2210.02156>.
- 438 [10] G. Christie, N. Fendley, J. Wilson, and R. Mukherjee. Functional map of the world, 2017. URL
439 <https://arxiv.org/abs/1711.07846>.
- 441 [11] A. Coates, A. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature
442 learning. In *Proceedings of the fourteenth international conference on artificial intelligence*
443 *and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings, 2011.
- 444 [12] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik. Emnist: an extension of mnist to handwritten
445 letters, 2017. URL <https://arxiv.org/abs/1702.05373>.
- 446 [13] F. Croce, M. Andriushchenko, V. Schwag, E. Debenedetti, N. Flammarion, M. Chiang, P. Mittal,
447 and M. Hein. Robustbench: a standardized adversarial robustness benchmark, 2021.
- 448 [14] A. Cutkosky and H. Mehta. Momentum improves normalized SGD. In H. D. III and A. Singh,
449 editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of
450 *Proceedings of Machine Learning Research*, pages 2260–2268. PMLR, 13–18 Jul 2020. URL
451 <https://proceedings.mlr.press/v119/cutkosky20b.html>.
- 452 [15] S. De, L. Berrada, J. Hayes, S. L. Smith, and B. Balle. Unlocking high-accuracy differentially pri-
453 vate image classification through scale, 2022. URL <https://arxiv.org/abs/2204.13650>.

- 454 [16] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical
455 image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages
456 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 457 [17] J. Diffenderfer, B. R. Bartoldson, S. Chaganti, J. Zhang, and B. Kailkhura. A winning hand:
458 Compressing deep networks can improve out-of-distribution robustness, 2021. URL <https://arxiv.org/abs/2106.09129>.
- 460 [18] J. Dong, A. Roth, and W. J. Su. Gaussian differential privacy, 2019. URL <https://arxiv.org/abs/1905.02383>.
- 462 [19] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Deghani,
463 M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16
464 words: Transformers for image recognition at scale, 2020. URL <https://arxiv.org/abs/2010.11929>.
- 466 [20] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private
467 data analysis. In *Proceedings of the Third Conference on Theory of Cryptography, TCC’06*,
468 page 265–284, Berlin, Heidelberg, 2006. Springer-Verlag. ISBN 3540327312.
- 469 [21] C. Dwork, N. Kohli, and D. Mulligan. Differential privacy in practice: Expose your epsilons!
470 *Journal of Privacy and Confidentiality*, 9, 10 2019. doi: 10.29012/jpc.689.
- 471 [22] U. Erlingsson, V. Feldman, I. Mironov, A. Raghunathan, K. Talwar, and A. Thakurta. Am-
472 plification by shuffling: From local to central differential privacy via anonymity, 2018. URL
473 <https://arxiv.org/abs/1811.12469>.
- 474 [23] H. Fang, X. Li, C. Fan, and P. Li. Improved convergence of differential private SGD with
475 gradient clipping. In *The Eleventh International Conference on Learning Representations, 2023*.
476 URL <https://openreview.net/forum?id=FRLswckPXQ5>.
- 477 [24] R. Geirhos, P. Rubisch, C. Michaelis, M. Bethge, F. A. Wichmann, and W. Brendel. Imagenet-
478 trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness,
479 2018. URL <https://arxiv.org/abs/1811.12231>.
- 480 [25] S. Ghalebikesabi, L. Berrada, S. Gowal, I. Ktena, R. Stanforth, J. Hayes, S. De, S. L. Smith,
481 O. Wiles, and B. Balle. Differentially private diffusion models generate useful synthetic images,
482 2023.
- 483 [26] A. Gohatkar, A. Achille, Y.-X. Wang, A. Roth, M. Kearns, and S. Soatto. Mixed differential
484 privacy in computer vision, 2022.
- 485 [27] S. Gopi, Y. T. Lee, and L. Wutschitz. Numerical composition of differential privacy, 2021. URL
486 <https://arxiv.org/abs/2106.02848>.
- 487 [28] P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia,
488 and K. He. Accurate, large minibatch sgd: Training imagenet in 1 hour, 2017. URL <https://arxiv.org/abs/1706.02677>.
- 490 [29] P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia,
491 and K. He. Accurate, large minibatch sgd: Training imagenet in 1 hour, 2018.
- 492 [30] S. Gupta, Y. Huang, Z. Zhong, T. Gao, K. Li, and D. Chen. Recovering private text in federated
493 learning of language models. In *Advances in Neural Information Processing Systems (NeurIPS)*,
494 2022.
- 495 [31] D. Hendrycks and T. Dietterich. Benchmarking neural network robustness to common corrup-
496 tions and perturbations, 2019. URL <https://arxiv.org/abs/1903.12261>.
- 497 [32] E. Hoffer, T. Ben-Nun, I. Hubara, N. Giladi, T. Hoefler, and D. Soudry. Augment your batch:
498 better training with larger batches, 2019. URL <https://arxiv.org/abs/1901.09335>.
- 499 [33] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora:
500 Low-rank adaptation of large language models, 2021.

- 501 [34] P. Izmailov, D. Podoprikin, T. Garipov, D. Vetrov, and A. G. Wilson. Averaging weights leads to
502 wider optima and better generalization, 2018. URL <https://arxiv.org/abs/1803.05407>.
- 503 [35] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford,
504 J. Wu, and D. Amodei. Scaling laws for neural language models, 2020. URL <https://arxiv.org/abs/2001.08361>.
- 506 [36] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization, 2014. URL <https://arxiv.org/abs/1412.6980>.
- 508 [37] B. Klimt and Y. Yang. The enron corpus: A new dataset for email classification research. In
509 *European conference on machine learning*, pages 217–226. Springer, 2004.
- 510 [38] P. W. Koh, S. Sagawa, H. Marklund, S. M. Xie, M. Zhang, A. Balsubramani, W. Hu, M. Ya-
511 sunaga, R. L. Phillips, I. Gao, T. Lee, E. David, I. Stavness, W. Guo, B. A. Earnshaw, I. S. Haque,
512 S. Beery, J. Leskovec, A. Kundaje, E. Pierson, S. Levine, C. Finn, and P. Liang. Wilds: A bench-
513 mark of in-the-wild distribution shifts, 2020. URL <https://arxiv.org/abs/2012.07421>.
- 514 [39] S. Kornblith, J. Shlens, and Q. V. Le. Do better imagenet models transfer better?, 2018. URL
515 <https://arxiv.org/abs/1805.08974>.
- 516 [40] A. Krizhevsky et al. Learning multiple layers of features from tiny images. Technical report,
517 Citeseer, 2009.
- 518 [41] A. Kumar, A. Raghunathan, R. Jones, T. Ma, and P. Liang. Fine-tuning can distort pretrained
519 features and underperform out-of-distribution, 2022. URL <https://arxiv.org/abs/2202.10054>.
- 520
- 521 [42] A. Kumar, A. Raghunathan, R. Jones, T. Ma, and P. Liang. Fine-tuning can distort pretrained
522 features and underperform out-of-distribution, 2022. URL <https://arxiv.org/abs/2202.10054>.
- 523
- 524 [43] A. Kumar, R. Shen, S. Bubeck, and S. Gunasekar. How to fine-tune vision models with sgd,
525 2022. URL <https://arxiv.org/abs/2211.09359>.
- 526 [44] N. Lambert, L. Castriato, L. von Werra, and A. Havrilla. Illustrating reinforcement learning
527 from human feedback (rlhf). *Hugging Face Blog*, 2022. <https://huggingface.co/blog/rlhf>.
- 528 [45] H. Li, P. Chaudhari, H. Yang, M. Lam, A. Ravichandran, R. Bhotika, and S. Soatto. Rethinking
529 the hyperparameters for fine-tuning, 2020. URL <https://arxiv.org/abs/2002.11770>.
- 530 [46] X. Li, F. Tramèr, P. Liang, and T. Hashimoto. Large language models can be strong differentially
531 private learners, 2021. URL <https://arxiv.org/abs/2110.05679>.
- 532 [47] X. Li, D. Liu, T. Hashimoto, H. A. Inan, J. Kulkarni, Y. T. Lee, and A. G. Thakurta. When does
533 differentially private learning not suffer in high dimensions?, 2022. URL <https://arxiv.org/abs/2207.00160>.
- 534
- 535 [48] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie. A convnet for the 2020s,
536 2022. URL <https://arxiv.org/abs/2201.03545>.
- 537 [49] X. Mao, Y. Chen, X. Jia, R. Zhang, H. Xue, and Z. Li. Context-aware robust fine-tuning, 2022.
538 URL <https://arxiv.org/abs/2211.16175>.
- 539 [50] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang. Learning differentially private recurrent
540 language models, 2017. URL <https://arxiv.org/abs/1710.06963>.
- 541 [51] H. Mehta, W. Krichene, A. Thakurta, A. Kurakin, and A. Cutkosky. Differentially private image
542 classification from features, 2022. URL <https://arxiv.org/abs/2211.13403>.
- 543 [52] H. Mehta, A. Thakurta, A. Kurakin, and A. Cutkosky. Large scale transfer learning for differen-
544 tially private image classification, 2022. URL <https://arxiv.org/abs/2205.02973>.

- 545 [53] R. Mehta, V. Albiero, L. Chen, I. Evtimov, T. Glaser, Z. Li, and T. Hassner. You only need
546 a good embeddings extractor to fix spurious correlations, 2022. URL [https://arxiv.org/
547 abs/2212.06254](https://arxiv.org/abs/2212.06254).
- 548 [54] S. Merity, C. Xiong, J. Bradbury, and R. Socher. Pointer sentinel mixture models. In *Inter-
549 national Conference on Learning Representations*, 2017. URL [https://openreview.net/
550 forum?id=Byj72udxe](https://openreview.net/forum?id=Byj72udxe).
- 551 [55] M. Nasr, S. Mahloujifar, X. Tang, P. Mittal, and A. Houmansadr. Effectively using public data
552 in privacy preserving machine learning, 2023. URL [https://openreview.net/forum?id=
553 5R96mIU85IW](https://openreview.net/forum?id=5R96mIU85IW).
- 554 [56] A. Panda, S. Mahloujifar, A. N. Bhagoji, S. Chakraborty, and P. Mittal. Sparsefed: Mitigating
555 model poisoning attacks in federated learning with sparsification, 2021. URL [https://arxiv.
556 org/abs/2112.06274](https://arxiv.org/abs/2112.06274).
- 557 [57] N. Papernot and T. Steinke. Hyperparameter tuning with renyi differential privacy, 2021. URL
558 <https://arxiv.org/abs/2110.03620>.
- 559 [58] Z. Peng, L. Dong, H. Bao, Q. Ye, and F. Wei. Beit v2: Masked image modeling with vector-
560 quantized visual tokenizers, 2022. URL <https://arxiv.org/abs/2208.06366>.
- 561 [59] B. Polyak and A. B. Juditsky. Acceleration of stochastic approximation by averaging. *Siam
562 Journal on Control and Optimization*, 30:838–855, 1992.
- 563 [60] V. Prabhu, S. Khare, D. Kartik, and J. Hoffman. Sentry: Selective entropy optimization via
564 committee consistency for unsupervised domain adaptation, 2020. URL [https://arxiv.org/
565 abs/2012.11460](https://arxiv.org/abs/2012.11460).
- 566 [61] N. Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12
567 (1):145–151, 1999.
- 568 [62] S. Qiao, H. Wang, C. Liu, W. Shen, and A. Yuille. Micro-batch training with batch-channel nor-
569 malization and weight standardization, 2019. URL <https://arxiv.org/abs/1903.10520>.
- 570 [63] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are
571 unsupervised multitask learners. 2019.
- 572 [64] H. Rahimian and S. Mehrotra. Frameworks and results in distributionally robust optimization.
573 *Open Journal of Mathematical Optimization*, 3:1–85, jul 2022. doi: 10.5802/ojmo.15. URL
574 <https://doi.org/10.5802/ojmo.15>.
- 575 [65] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis
576 with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision
577 and Pattern Recognition (CVPR)*, pages 10684–10695, June 2022.
- 578 [66] E. K. Ryu and S. P. Boyd. A primer on monotone operator methods. 2015.
- 579 [67] S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang. Distributionally robust neural networks
580 for group shifts: On the importance of regularization for worst-case generalization, 2019. URL
581 <https://arxiv.org/abs/1911.08731>.
- 582 [68] V. Shejwalkar, A. Ganesh, R. Mathews, O. Thakkar, and A. Thakurta. Recycling scraps:
583 Improving private learning by leveraging intermediate checkpoints, 2022. URL [https://
584 arxiv.org/abs/2210.01864](https://arxiv.org/abs/2210.01864).
- 585 [69] Y. Shen, Z. Wang, R. Sun, and X. Shen. Towards understanding the impact of model size on
586 differential private classification, 2021. URL <https://arxiv.org/abs/2111.13895>.
- 587 [70] W. Shi, S. Chen, C. Zhang, R. Jia, and Z. Yu. Just fine-tune twice: Selective differential privacy
588 for large language models. *arXiv preprint arXiv:2204.07667*, 2022.
- 589 [71] C. Shorten and T. M. Khoshgoftaar. A survey on image data augmentation for deep learning.
590 *Journal of Big Data*, 6(1):60, Jul 2019. ISSN 2196-1115. doi: 10.1186/s40537-019-0197-0.
591 URL <https://doi.org/10.1186/s40537-019-0197-0>.

- 592 [72] S. Song, K. Chaudhuri, and A. D. Sarwate. Stochastic gradient descent with differentially
593 private updates. In *2013 IEEE Global Conference on Signal and Information Processing*, pages
594 245–248, 2013. doi: 10.1109/GlobalSIP.2013.6736861.
- 595 [73] I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and
596 momentum in deep learning. In S. Dasgupta and D. McAllester, editors, *Proceedings of the*
597 *30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine*
598 *Learning Research*, pages 1139–1147, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL
599 <https://proceedings.mlr.press/v28/sutskever13.html>.
- 600 [74] A. Team. Learning with privacy at scale, 2017. URL [https://docs-assets.developer.](https://docs-assets.developer.apple.com/ml-research/papers/learning-with-privacy-at-scale.pdf)
601 [apple.com/ml-research/papers/learning-with-privacy-at-scale.pdf](https://docs-assets.developer.apple.com/ml-research/papers/learning-with-privacy-at-scale.pdf).
- 602 [75] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jegou. Training data-
603 efficient image transformers; distillation through attention. In M. Meila and T. Zhang, editors,
604 *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Pro-*
605 *ceedings of Machine Learning Research*, pages 10347–10357. PMLR, 18–24 Jul 2021. URL
606 <https://proceedings.mlr.press/v139/touvron21a.html>.
- 607 [76] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman. GLUE: A multi-task bench-
608 mark and analysis platform for natural language understanding. In *International Conference on*
609 *Learning Representations*, 2019. URL <https://openreview.net/forum?id=rJ4km2R5t7>.
- 610 [77] Y.-X. Wang, B. Balle, and S. P. Kasiviswanathan. Subsampled renyi differential privacy and
611 analytical moments accountant. In K. Chaudhuri and M. Sugiyama, editors, *Proceedings of the*
612 *Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of
613 *Proceedings of Machine Learning Research*, pages 1226–1235. PMLR, 16–18 Apr 2019. URL
614 <https://proceedings.mlr.press/v89/wang19b.html>.
- 615 [78] R. Wightman. Pytorch image models. [https://github.com/rwightman/](https://github.com/rwightman/pytorch-image-models)
616 [pytorch-image-models](https://github.com/rwightman/pytorch-image-models), 2019.
- 617 [79] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf,
618 M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. L. Scao,
619 S. Gugger, M. Drame, Q. Lhoest, and A. M. Rush. Huggingface’s transformers: State-of-the-art
620 natural language processing, 2019. URL <https://arxiv.org/abs/1910.03771>.
- 621 [80] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking
622 machine learning algorithms, 2017. URL <https://arxiv.org/abs/1708.07747>.
- 623 [81] A. Yousefpour, I. Shilov, A. Sablayrolles, D. Testuggine, K. Prasad, M. Malek, J. Nguyen,
624 S. Ghosh, A. Bharadwaj, J. Zhao, G. Cormode, and I. Mironov. Opacus: User-friendly
625 differential privacy library in pytorch, 2021. URL <https://arxiv.org/abs/2109.12298>.
- 626 [82] D. Yu, S. Naik, A. Backurs, S. Gopi, H. A. Inan, G. Kamath, J. Kulkarni, Y. T. Lee, A. Manoel,
627 L. Wutschitz, S. Yekhanin, and H. Zhang. Differentially private fine-tuning of language models,
628 2021. URL <https://arxiv.org/abs/2110.06500>.
- 629 [83] D. Yu, H. Zhang, W. Chen, J. Yin, and T.-Y. Liu. Large scale private learning via low-rank
630 reparametrization, 2021.
- 631 [84] S. Zhang, E. Dinan, J. Urbanek, A. Szlam, D. Kiela, and J. Weston. Personalizing dialogue
632 agents: I have a dog, do you have pets too?, 2018.
- 633 [85] S. Zhang, E. Dinan, J. Urbanek, A. Szlam, D. Kiela, and J. Weston. Personalizing dialogue
634 agents: I have a dog, do you have pets too?, 2018. URL [https://arxiv.org/abs/1801.](https://arxiv.org/abs/1801.07243)
635 [07243](https://arxiv.org/abs/1801.07243).
- 636 [86] Y. Zhu and Y.-X. Wang. Poission subsampled rényi differential privacy. In K. Chaudhuri
637 and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine*
638 *Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7634–7642. PMLR,
639 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/zhu19c.html>.