
Label Privacy in Split Learning for Large Models with Parameter-Efficient Training

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 As deep learning models become larger and more expensive, many practitioners
2 turn to fine-tuning APIs. These web services allow fine-tuning a model between
3 two parties: the client that provides the data, and the server that hosts the model.
4 While convenient, these APIs raise a new concern: the data of the client is at
5 risk of privacy breach during the training procedure. This challenge presents
6 an important practical case of vertical federated learning, where the two parties
7 perform parameter-efficient fine-tuning (PEFT) of a large model. In this study, we
8 systematically search for a way to fine-tune models over an API *while keeping the*
9 *labels private*. We analyze the privacy of LoRA, a popular approach for parameter-
10 efficient fine-tuning when training over an API. Using this analysis, we propose
11 P³EFT, a multi-party split learning algorithm that takes advantage of existing PEFT
12 properties to maintain privacy at a lower performance overhead. To validate our
13 algorithm, we fine-tune DeBERTa-v2-XXLarge, Flan-T5 Large and LLaMA-2 7B
14 using LoRA adapters on a range of NLP tasks. We find that P³EFT is competitive
15 with existing privacy-preserving methods in multi-party and two-party setups while
16 having higher accuracy.

17 1 Introduction

18 One of the main reasons behind deep learning success is its ability to transfer knowledge between
19 tasks [34]. When training a model for any particular problem, it is common to reuse previously
20 trained models from other, related problems. In the past, this was typically done by downloading
21 pre-trained model weights from public hubs, then fine-tuning the said models on the downstream
22 task. However, as models grow larger and more compute-intensive, fine-tuning them locally becomes
23 an increasingly difficult task. Furthermore, many recent models are not released, but instead made
24 available as proprietary services.

25 When a model cannot be fine-tuned locally, many practitioners opt instead for the so-called fine-tuning
26 APIs [27, 16, 6, 26]. These APIs are web services that host one or several pre-trained models and
27 allow clients to perform limited fine-tuning. More specifically, APIs usually allow their clients to run
28 parameter-efficient fine-tuning (PEFT), such as LoRA [15] or Prefix-tuning [21]. These techniques
29 allow adapting a model to a dataset while training a relatively small number of additional weights,
30 which is particularly important for large language or image generation models that have billions of
31 parameters.

32 Although the fine-tuning APIs can be convenient, they also introduce new risk in terms of data privacy.
33 When a client uses such API to train on sensitive data, they need to ensure that their data will stay
34 private [7]. This is particularly important when dealing with patient’s medical records, personal
35 user data or trade secrets [24, 19]. The two main threats to data privacy are that the API provider
36 obtains the private data and that a third party intercepts data in transit. Therefore, data privacy is

37 not guaranteed even if the API provider is trusted. Several recent works propose LLM fine-tuning
38 protocols that establish a certain level of privacy for multi-party fine-tuning [42, 7, 22]. Unfortunately,
39 these algorithms work for a narrow class of fine-tuning algorithms or assume that a client can run
40 LLM training locally using an obfuscated version of the model, provided by a remote server [42].
41 As a result, these algorithms are impractical for our use case of fine-tuning over an API. The few
42 algorithms that are suitable for API fine-tuning guarantee the privacy of input tokens [22], meaning
43 that the attacker can infer private training *labels*.

44 In this work, we seek to alleviate this problem by designing a two-party fine-tuning protocol that
45 performs standard parameter-efficient fine-tuning with privacy guarantees. We formulate our protocol
46 as a **special case of split learning** (or vertical federated learning), where one side (server) holds the
47 pre-trained model and the other (client) has private training data. More specifically, we focus on **the**
48 **privacy of client’s training labels**. While input privacy is often crucial, there are scenarios where
49 input data is publicly available, such as social media user pages. In these cases, labels could include
50 ad clicks (known to the social network) or financial information (known to a bank that matches social
51 profiles to its internal records). This example further justifies the use of LLMs, as social media pages
52 often contain substantial amounts of text, and LLMs excel at processing long-context data.

53 Instead of developing a specific privacy-preserving architecture, we seek algorithms that can work
54 with popular existing models and PEFT algorithms. Furthermore, our approach relies on the properties
55 of parameter-efficient fine-tuning. Notably, since the adapters are compact, both parties can maintain
56 multiple sets of adapters and swap between them with relative ease. This allows us to design a
57 PEFT-specific algorithm that can solve its task more effectively than general split learning strategies
58 [18].

59 We summarize our main contributions as follows:

- 60 • We analyze Low-Rank Adaptation, a common parameter-efficient fine-tuning algorithm,
61 from the perspective of label privacy in the split learning setup. We observe that, despite
62 fine-tuning less than 0.1% of model parameters, PEFT algorithms leak client’s training
63 labels against simple attacks that work for modern pretrained transformers.
- 64 • Based on our analysis, we formulate a framework for privacy-preserving parameter-efficient
65 fine-tuning (P³EFT). This framework leverages the properties of PEFT to obfuscate the
66 gradients and parameters communicated during fine-tuning with little impact on the fine-
67 tuned model quality.
- 68 • To verify the practical viability of P³EFT, we conduct experiments on popular real-world
69 PEFT workloads¹. Specifically, we fine-tune DeBERTa-v2-XXL [13], Flan-T5-Large [4]
70 and LLaMA-2 7B [35] on a set of standard language understanding problems. We find that,
71 compared to prior split learning algorithms, P³EFT can maintain label privacy throughout
72 training with a significantly smaller accuracy drop.

73 2 Background

74 2.1 Federated learning and split learning

75 Privacy preservation in machine learning has been a subject of active study within several frameworks.
76 An important branch of privacy-preserving learning methods is federated learning, or FL [24], which
77 can be broadly described as an approach allowing several parties to train a model jointly without
78 sharing their private data. In particular, vertical federated learning [12, 43] targets the scenario where
79 different features (including the label) of each training instance are kept by different parties.

80 One of the most popular approaches to vertical FL for neural networks is split learning [10, 37],
81 where each party stores its part of the overall model. To train the model in such an approach, it is
82 only necessary to transfer the intermediate activations and the gradients between layers, while the
83 data itself is stored at the premises of the participant hosting each layer. In this work, we focus on
84 the two-party formulation of split learning, where one side stores the features for each example and
85 another one stores the labels.

¹The code is available at github.com/anonymousauthor56/P3EFT

86 Recent works have investigated the setting of two-party split learning from the label leakage per-
87 spective [38, 28]: because the label party needs to pass the gradients of the loss function to the
88 non-label party, it is possible for the latter party to deduce the labels by inspecting the gradients or
89 activations or by hijacking the training procedure. Li et al. [18] provide a set of attack methods that
90 allow recovering private labels and propose a defense mechanism that injects noise into the gradients;
91 however, they test the approach on pretraining smaller models, and we study finetuning large models
92 on private downstream data.

93 2.2 Parameter-efficient finetuning

94 The majority of large neural networks today are not trained with a specific task in mind: instead, they
95 are pretrained on a general objective and then adapted for the downstream problem. Importantly, the
96 growth in the size of foundation models has led to the increased popularity of parameter-efficient
97 finetuning (PEFT) methods that adapt the model to a given task by training a small number of
98 task-specific parameters. There are several prominent approaches to parameter-efficient finetuning,
99 ranging from trainable prompts [21, 11], to residual adapters [14, 29]. We focus on Low-Rank
100 Adaptation (or LoRA, [15]), one of the most popular PEFT methods that adds extra parameters to each
101 weight matrix in the form of a low-rank factorization (see Appendix C for a more detailed description).
102 Such formulation allows LoRA adapters to be merged into the original weights after finetuning; this
103 ability, combined with the simplicity of the method, has made LoRA a broadly popular approach in
104 multiple domains. Still, the approach we propose can be applied to any PEFT method.

105 Several recent lines of work explore the problem of fine-tuning LLMs with privacy guarantees [44, 31].
106 Zhao et al. [46] analyze the viability of prompt tuning for federated learning, and Zhang et al. [45], Liu
107 et al. [23] study PEFT algorithms in the setting of *horizontal* federated learning, that is, where multiple
108 users train a shared model on their local private data. Another, more relevant research direction
109 considers private fine-tuning in a *vertical* federated learning scenario, where participants hold different
110 model layers [22, 40]. Most of these studies leverage the idea of differential privacy to prove an
111 upper bound on how much information is leaked [8]. Unfortunately, these upper bounds are typically
112 loose and do not match practical observations for real models. Furthermore, the majority of these
113 studies only guarantees privacy of specific parts of the training procedure: for instance, Li et al.
114 [22] only protects the input features, and not labels or model parameters. Finally, Xiao et al. [42]
115 presents an alternative algorithm that protects client data by running the entire fine-tuning on client
116 side by emulating the server-side model layers. While this approach is more holistic, it assumes that
117 clients can run fine-tuning locally, which makes it impractical for many real-world users of LLM
118 fine-tuning APIs. The primary distinction between our work and these studies is that we investigate
119 parameter-efficient adaptation in the setting of split learning: we aim to finetune a model without
120 disclosing the labels of examples to the model provider.

121 3 Privacy-preserving parameter-efficient fine-tuning

122 In this section, we analyze the privacy of parameter-efficient fine-tuning and propose a protocol for
123 two-party parameter-efficient fine-tuning with the desired privacy guarantees. We begin by analyzing
124 the privacy of API fine-tuning with popular PEFT algorithms in Sections 3.1 and 3.2. Then, in
125 Section 3.3, we formulate a protocol for privately computing gradients over fine-tuning APIs. Finally,
126 we formulate the full P³EFT protocol in Section 3.4.

127 3.1 Setup

128 To analyze the privacy of API fine-tuning, we first need to formulate a common framework for
129 this type of APIs and develop private learning protocols. This step is important, because existing
130 fine-tuning APIs greatly vary in what they offer to the client: from closed APIs that require users to
131 submit their full training data [27] to more flexible APIs where clients can run individual training
132 steps [20, 2, 30]. Similarly to most existing works on split learning, we focus on the latter type of
133 APIs that allows clients to run individual forward and backward passes over a remote model. Thus,
134 a client can use these APIs to obtain the training gradients for their PEFT adapters, then update
135 adapters locally with any optimization method. In our work, we adopt this archetype of fine-tuning
136 API as it offers sufficient flexibility to develop privacy-preserving algorithms.

137 We formulate fine-tuning over an API for two or more parties: a client, and one or several servers.
138 The client owns a training dataset with inputs X and labels Y . In turn, each server has the same
139 pre-trained model $h(x_i, \theta) \in \mathcal{R}^d$. Note that the parameters θ denote not the pre-trained model

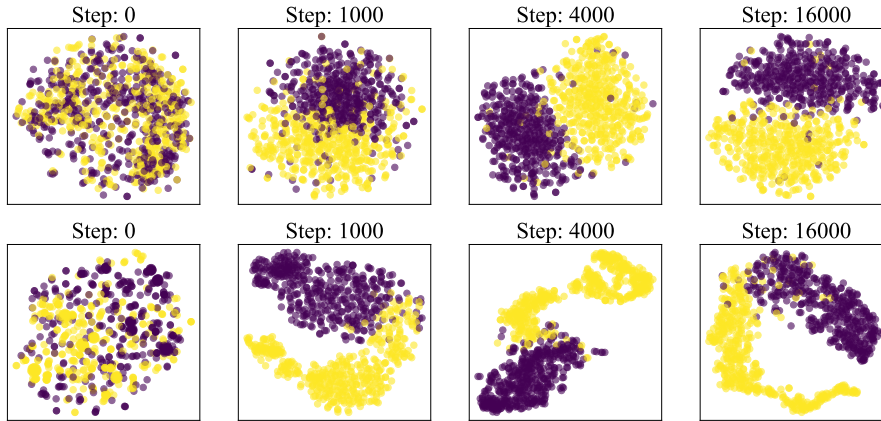


Figure 1: A visualization of top-2 principal components of gradients (top) and activations (bottom) from different fine-tuning steps (left to right). Color indicates the training labels (binary).

140 weights, but the trainable adapter weights for a certain PEFT algorithm. A model can encode an input
 141 $x_i \in X$ and produce a d -dimensional vector of activations that depend on the learned adapter weights
 142 θ .

143 To allow fine-tuning, a server offers two API methods:

- 144 1. **forward** $(x, \theta) \rightarrow h(x, \theta)$ that computes model activations on input x using adapter weights
 145 θ ;
- 146 2. **backprop** $(x, \theta, g_h) \rightarrow g_\theta$ that receives gradients of an arbitrary loss function w.r.t. model
 147 activations $g_h = \frac{\partial L(h(x, \theta))}{\partial h(x, \theta)}$ and returns the gradients w.r.t. adapter parameters, $g_\theta =$
 148 $\frac{\partial L(h(x, \theta))}{\partial \theta}$.

149 We further assume that both `forward`(\cdot) and `backprop`(\cdot) APIs are stateless and deterministic, i.e.
 150 calling the same API method multiple times (or on multiple servers) with the same inputs produces
 151 identical results. Thus, if the model uses dropout or any other form of non-determinism, we assume
 152 that clients provide the random seed as a part of x .

153 To fine-tune a model with this API, a client can initialize adapters locally, alongside with a small
 154 task-specific head², then train both adapters and the head. For each training batch $(x, y) \in D$, a client
 155 calls `forward` (x, θ) to compute feature representations, then predicts with local “head” and computes
 156 task-specific loss function L . After that, a client performs backward pass: first, it computes gradients
 157 w.r.t. local head inputs $g_h = \frac{\partial L}{\partial h}$, then passes those gradients to a remote server via `backprop` (x, θ, g_h)
 158 API call to compute gradients w.r.t. $\frac{\partial L}{\partial \theta}$. Finally, a client updates both θ and local “head” parameters
 159 using the optimizer of choice.

160 Before building more advanced algorithms, let us analyze the privacy of client’s labels under standard
 161 fine-tuning. We consider an “honest, but curious” attacker model. This means that the server will
 162 faithfully run the forward and backprop computations as requested by the client without changing
 163 the results. Furthermore, we assume that servers are independent and do not communicate client’s
 164 data between each other. However, a server can recover client’s labels by performing arbitrary
 165 computations using any information it receives from the client. When training in this way, a client
 166 does not directly communicate training labels to the server. However, it communicates inputs, adapter
 167 parameters, and gradients. Furthermore, the server communicates input representations that can be
 168 intercepted by a third party.

169 3.2 Label Leakage of Standard Split Learning

170 In Figure 1, we train a DeBERTa-v2-XXL model on the SST-2 [32] sentiment classification dataset.
 171 The top row depicts the gradients g_h communicated by the client when calling `backprop`(\cdot) at different
 172 training stages. In the bottom row, we similarly track activations $h(x, \theta)$ that server may compute
 173 based on the specified x, θ . We defer further additional figures and details to Section 4.1.

174 As we can see, both gradients and activations are arranged in such a way that simple k-means
 175 clustering would reveal which objects have the same label. The training activations (bottom row) do

²A linear layer that predicts class logits or regression target.

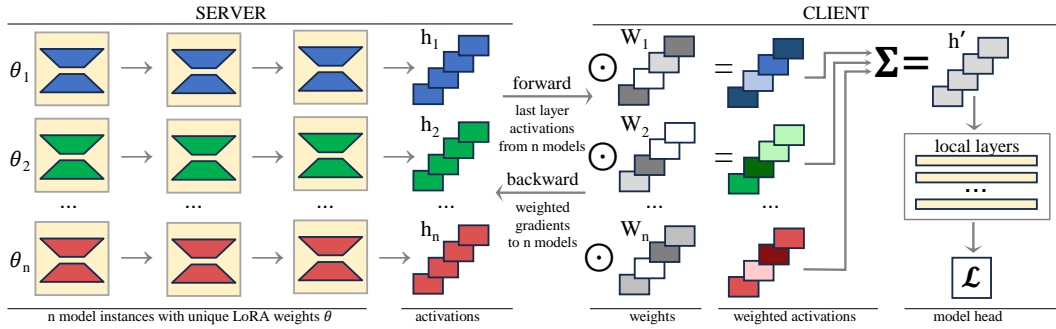


Figure 2: An intuitive illustration of the proposed fine-tuning protocol.

176 not reveal labels right away (at least not against this attack). However, they gradually “leak” private
 177 label information during training. Informally, it appears that the training gradients gradually pull
 178 apart the feature representations for each label, until eventually they turn into separate clusters. From
 179 an information-theoretic perspective, knowing just one vector of gradients *or* trained activations
 180 allows the attacker to learn all but one bit³ of information about client’s private labels.

181 To summarize, leaving any *one* data source unprotected (gradients, activations or parameters) would
 182 already compromise label privacy. However, we found that gradients and activations require different
 183 means of protection.

184 3.3 Privacy-preserving backpropagation

185 In this section, we formulate an algorithm for “anonymizing” the gradients communicated over a
 186 single training step with arbitrary PEFT type. Several prior works approach this by modifying the
 187 training objective or model architecture. However, when dealing with a real-world PEFT workload
 188 with optimized hyperparameters, changing the model or loss function often results in reduced model
 189 accuracy⁴. Thus, we seek an algorithm that preserves both model and training objective.

190 We design our algorithm based on an observation that **backpropagation is conditionally lin-**
 191 **ear in output gradients**, even when the model itself is nonlinear. Formally, if we take a model
 192 $h(\cdot, \cdot)$, a fixed set of trainable parameters θ and input samples x , the backprop function⁵ computes
 193 $\text{backprop}(x, \theta, \frac{\partial L}{\partial h(x, \theta)}) = \frac{\partial L}{\partial \theta}$. For convenience, we shorten it to $\text{backprop}(x, \theta, g_h) = g_\theta$, where
 194 $g_h = \frac{\partial L}{\partial h(x, \theta)}$ represents the gradients of some objective function with respect to model activations
 195 (outputs), and $g_\theta = \frac{\partial L}{\partial \theta}$ are gradients of the same objective function w.r.t. trainable parameters. In
 196 this notation, backprop is linear in terms of g_h for any fixed x, θ .

197 This becomes self-evident if we view backprop as multiplying g_h by the Jacobian of model outputs
 198 w.r.t. trainable parameters, $\frac{\partial h(x, \theta)}{\partial \theta}$. If x, θ are constant, the Jacobian is also constant, and backprop
 199 is a linear operator:

$$\text{backprop}(x, \theta, \frac{\partial L}{\partial h(x, \theta)}) = \frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial h(x, \theta)} \times \frac{\partial h(x, \theta)}{\partial \theta}. \quad (1)$$

200 This observation allows us to design a private backpropagation protocol. To illustrate
 201 this protocol, let us first consider a distributed API with two identical independent servers
 202 that offer backprop API. Then, for arbitrary vector z , we can rewrite $\text{backprop}(x, \theta, g_h)$ as
 203 $\text{backprop}(x, \theta, g_h + z) + \text{backprop}(x, \theta, g_h - z)$.

204 During API fine-tuning, we obtain $\text{backprop}(x, \theta, g_h + z)$ using an API call to server 1, whereas the
 205 second term $\text{backprop}(x, \theta, g_h - z)$ translates to an API call to server 2. Note that neither of two
 206 servers has access to the true gradient g_h : they only receive the sum $[g_h + z]$. If we sample a large
 207 noise vector z ($\text{Var}(z) \gg \|g_h\|_2^2$), this sum also becomes dominated by noise. However, when both
 208 API calls finish, a client can sum the results to recover the true gradient of the loss with respect to
 209 parameters.

210 If both requests are processed by the same server, it can obviously recover g_h by adding up gradients
 211 from both calls, which leads us to the final step. Instead of generating a single noise vector, a client

³The missing bit corresponds to attacker not knowing which cluster corresponds to label “1”.

⁴We validate this empirically in 4.2.

⁵This is the same as the backprop API defined in Section 3.1.

212 needs to generate (privately) a set of $m > 1$ random vectors $\hat{g}_h^1, \dots, \hat{g}_h^m$ and scalars $\alpha_1, \dots, \alpha_m$ such
 213 that

$$g_h = \sum_{i=1}^m \alpha_i \cdot \hat{g}_h^i. \quad (2)$$

214 Then, for each \hat{g}_h^i , client computes $\text{backprop}(x, \theta, \hat{g}_h^i)$ as m parallel API calls. Once this is done,
 215 client recovers

$$g_\theta = \sum_{i=1}^m \alpha_i \cdot \text{backprop}(x, \theta, \hat{g}_h^i). \quad (3)$$

216 Note that the client does not reveal $\alpha_1, \dots, \alpha_m$ to anyone.

217 The resulting procedure is formulated in Algorithm 1. This algorithm is conceptually similar to
 218 the secure aggregation protocol for conventional (horizontal) federated learning [1]. This protocol
 219 allows clients to average their local vector with peers while keeping each individual vector provably
 220 private. Similarly to our scheme, clients perturb the vector in such a way that the average of perturbed
 221 vectors remains the same. Unlike Bonawitz et al. [1], our protocol privately backpropagates through
 222 a server-hosted model by leveraging the conditional linearity of the backpropagation operator.

Algorithm 1 private_backprop — Privacy-Preserving Backpropagation (from the client’s perspective)

```

1: Input:  $x$  inputs,  $\theta$  adapter weights,  $g_h$  gradients w.r.t. activations,  $m > 1$  - number of passes
2:  $\hat{g}_h^1, \dots, \hat{g}_h^m, \alpha_1, \dots, \alpha_m = \text{obfuscate}(g_h, m)$  ▷ 2
3: for  $j = 1, \dots, m$  do
4:    $\hat{g}_\theta^j = \text{backprop}(x, \theta, \hat{g}_h^j)$  ▷ computed by server
5: end for
6:  $g_\theta = \sum_{j=1}^m \alpha_j \cdot \hat{g}_\theta^j$ 
7: Return:  $g_\theta$ 

```

223 The private backpropagation algorithm can allow client to safely compute gradients *once*, but, in
 224 practice, client usually needs to run many consecutive steps. This creates an additional vector of
 225 attack: if the same server receives two sets of parameters θ_t, θ_{t+1} , they could potentially recover g_θ
 226 by inverting the optimizer.

227 In the simplest case, if the server somehow knows that the client computes $\theta_{t+1} = \theta_t - \eta \cdot g_\theta$, then
 228 they can compute $g_\theta = (\theta_t - \theta_{t+1})/\eta$. While g_θ does not necessarily leak private labels, a server
 229 could, in some cases, use g_θ to recover g_h , either fully (e.g. if Jacobian is invertible), or partially.

230 The client has two ways to prevent this attack. The first one is to ensure that no single server runs
 231 backprop on two consecutive steps. This is easy to do in decentralized systems where there are
 232 many potential servers. However, even when there is a single server, they could be required to set up
 233 multiple trusted execution environments [25]. A more risky alternative is to ensure that the gradients
 234 cannot be reversed from consecutive parameters: randomize initial optimizer statistics or add noise to
 235 parameters. This solution is easier, but it can slow down training in some cases.

236 To summarize, we formulated a procedure that allows a client to compute gradients privately for any
 237 given model and PEFT type. Furthermore, since Equation 3 recovers true gradients, this obfuscation
 238 method does not affect the training dynamics. However, as we have shown in Section 3.1, gradients
 239 are not the only source of privacy leakage.

240 3.4 Full fine-tuning

241 The other major attack vector are training activations. As the model fits to training data, it’s
 242 intermediate activations $h(x, \theta)$ allow attackers to recover labels, e.g. by clustering (see Figure 1).
 243 To combat this issue, we take advantage of the fact that PEFT has few trainable parameters. Instead
 244 of learning just one set of trainable parameters, a client creates n independent adapter sets $\theta_1, \dots, \theta_n$.
 245 Note that this does not require n unique servers: a single server can run multiple sets of adapters.
 246 Furthermore, a client can alternate between using different servers for the same adapters. During
 247 forward pass, the outputs of different adapters are mixed together using randomized mixing weights
 248 $W \in \mathcal{R}^{n,d}$:

$$h'(x, \theta_1, \dots, \theta_n) = \sum_{i=1}^n W_i \odot h(x, \theta_i) \quad (4)$$

249 Overall, we design this model in such a way the combined model h' can predict the labels, but the
 250 adapters $h(x, \theta_i)$ do not allow predicting these labels without knowing the mixing weights W . The
 251 mixing weights are generated such that initial activations $h'(x, \dots)$ are equal to mean $h(x, \cdot)$ for all
 252 x . To achieve this, we generate W as follows: first, we generate $n \cdot (n - 1)/2$ d-dimensional random
 253 vectors $\xi_{i,j} \in \mathcal{R}^d \forall i \in [1, n], j \in [i + 1, n]$. Then, we add them up in the following way:
 254

$$W = \begin{pmatrix} \frac{1}{n}e + \xi_{1,2} + \xi_{1,3} + \dots + \xi_{1,n} \\ -\xi_{1,2} + \frac{1}{n}e + \xi_{2,3} + \dots + \xi_{2,n} \\ \dots \\ -\xi_{1,n} - \xi_{2,n} - \xi_{3,n} - \dots + \frac{1}{n}e \end{pmatrix} \quad (5)$$

255 Here, e stands for a vector of all ones. The purpose of these mixing weights is to ensure that the
 256 gradients w.r.t. individual $h(x, \theta_i)$ are obfuscated, but the averaged model behaves the same as
 257 regular PEFT adapter. To illustrate this, consider $n=2$ identical LoRA adapters θ_1, θ_2 . During the
 258 first training step $h(x, \theta_1) = h(x, \theta_2)$. Therefore,

$$h'(x, \theta_1, \dots, \theta_n) = (1/2e + \xi_{1,2}) \odot h(x, \theta_1) + (1/2e - \xi_{1,2}) \odot h(x, \theta_2) = h(x, \theta_1) \quad (6)$$

259 However, the two adapters will learn different functions as they receive different gradients. From the
 260 first update on, h' will be equal to an average of adapter predictions.

261 Finally, to ensure that individual adapters $h(x, \theta)$ do not accidentally “learn to leak” labels, we
 262 maintain this over the course of training with a privacy regularizer inspired by [9]. This ensures that
 263 it is impossible to predict labels from individual adapters $h(x, \theta_i)$. Intuitively, on each training step,
 264 client fits n linear “heads” that learn to predict labels y from $h(x, \theta_i)$, then performs an adversarial
 265 update of θ_i to prevent the “head” from predicting y . Formally, each of n “heads” minimize the same
 266 objective function as the full model. For instance, if the full model solves multi-class classification,
 267 each head is trained to minimize cross-entropy:

$$\eta_i^* = \arg \min_{\eta_i} \sum_{x,y \in D} -y \cdot \log \frac{e^{\langle \eta_{ij}, h(x, \theta_i) \rangle}}{\sum_k e^{\langle \eta_{ik}, h(x, \theta_i) \rangle}}, \quad (7)$$

268 where y is one-hot encoding of the correct class.
 269

270 The whole adversarial update takes place locally on client’s side, using the same $h(x, \theta)$ it uses for the
 271 main training objective. The resulting procedure appears complicated but it typically takes negligible
 272 time compared to running the large pre-trained model $h(x, \theta)$. Furthermore, since adversarial “heads”
 273 are linear, minimizing the objective above is done with standard logistic regression solver.

274 To summarize, our approach combines the two proposed ideas: we use the private backpropagation
 275 algorithm from Section 3.3 to protect the gradients, then trains a mixture of adapters in such a way
 276 that obfuscates learned activations leaking labels. The resulting procedure is described in Algorithm 2.
 277 In the next section, we will evaluate the efficacy of P³EFT on popular NLP benchmarks.

278 4 Experiments

279 The main goal of our study is to find a practical method of private fine-tuning that would scale to
 280 large models. Because our approach leverages parameter-efficient fine-tuning techniques, we evaluate
 281 P³EFT with fine-tuning Transformer models on popular NLP benchmarks that these techniques were
 282 designed for.

283 To that end, we chose three pre-trained models: DeBERTa-XXLarge [13], Flan-T5-Large [4] and
 284 LLaMA-2 7B [35]. We train these models on several datasets from the GLUE benchmark [39]:
 285 SST-2 [32], MNLI [41] and QNLI.

286 4.1 Privacy of gradients and activations

287 For this experiment, we train DeBERTa-XXLarge on SST-2 dataset using LoRA adapters with
 288 hyperparameters from [15]. First, we train the model locally and track model activations h and
 289 gradients w.r.t. those activations. We apply principal component analysis to them and plot the first

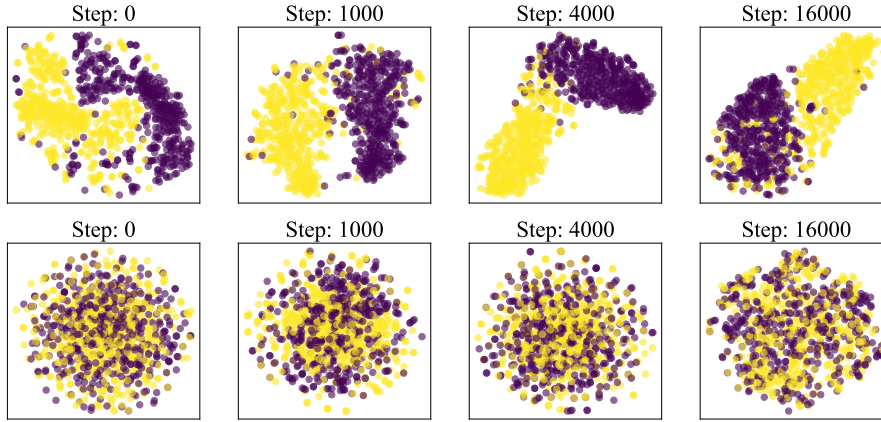


Figure 3: Gradients of cross-entropy w.r.t. LoRA parameters for DeBERTa-v2-XXLarge. The top row corresponds to normal backpropagation and the bottom row uses privacy-preserving backprop.

290 2 dimensions in Figure 1. Similarly, we visualize gradients of individual per-sample loss functions
 291 w.r.t. LoRA parameters θ in Figure 3 (top row). The results suggest that a hypothetical attacker could
 292 easily recover private labels by performing K-Means clustering over any data source: activations,
 293 gradients with respect to activations, or individual gradients with respect to parameters.

294 Next, we run the same experiment using privacy-preserving backpropagation as defined in Section 3.3.
 295 We use $n = 2$ with the noise variance set to 1000. As expected, we observed the same learning curve
 296 as with normal training. However, instead of sending gradients w.r.t. activations to the server, a
 297 client uses specially crafted random noise vectors that are not informative. In Figure 3 (bottom) we
 298 plot the same kind of individual gradients as in the top row, except that we visualize the gradients
 299 computed by the first of the two servers. Finally, we train XGBoost [3] with default hyperparameters
 300 to predict labels given the noisy gradients (pre-PCA): the resulting classifier is able to fit the training
 301 data perfectly, but has at most 50.4% accuracy on a balanced test set.

302 4.2 Main fine-tuning experiments

303 Next, we evaluated the entire P3EFT algorithm. To control tasks and model type, we examined
 304 DeBERTa and Flan-T5 across all four datasets mentioned above, in addition to evaluating LLaMA on
 305 SST2 and QNLI datasets. For each setup, we compare against three baselines:

- 306 • **Without LoRAs.** In this baseline, the client gathers h activations at the beginning (with no
 307 adapters), then proceeds to train local “head” layers using these activations. This method cannot
 308 leak information about training labels except for what is stored in X .
- 309 • **Regular fine-tuning (Regular FT)** refers to training a single LoRA adapter normally. This baseline
 310 represents an upper bound on model accuracy, but lacks privacy.
- 311 • **Distance Correlation (DC).** Our re-implementation of the distance correlation defense formulated
 312 in [33] for Transformer models.

313 For each algorithm, we evaluated a task-specific metric (accuracy or F1), as well as the privacy
 314 leakage value for the 3 following measures:

- 315 • **Spectral attack AUC** — a measure of vulnerability to an attack proposed in [33], measured as
 316 classifier ROC AUC: lower value corresponds to better privacy.
- 317 • **Norm attack AUC** — vulnerability to a variant of attack proposed in [18], measured as classifier
 318 ROC AUC (lower is better). Despite the initial proposal of this approach for attacking gradients,
 319 we observed that it is also well-suited for attacking activations.
- 320 • **K-means accuracy** — vulnerability to clusterization attack, measured in the percentage of correctly
 321 clustered activations, lower is better.

322 For all setups, we report the worst (least private) value among these metrics throughout the entire
 323 training period as a measure of privacy leakage, because it is the worst possible scenario that matters
 324 from the client’s perspective. For DC and P³EFT, we specify the values for the best configuration in
 325 terms of the utility-privacy trade-off. See details in Appendix A. We also report adjusted standard
 326 deviations for the two privacy aware algorithms: P³EFT and DC. To do so, we run the full training
 327 procedure from scratch with 3 random seeds.

Table 1: Accuracy and privacy metrics.
DeBERTa XXLLarge.

Dataset		Without LoRAs	Regular FT	DC	P ³ EFT
SST2	<i>acc</i>	82.9	96.9	96.6 \pm 0.4	96.5 \pm 0.2
	<i>leak</i>	53.9	99.1	93.3 \pm 6.8	62.6 \pm 2.6
QNLI	<i>acc</i>	72.6	96.0	95.8 \pm 0.3	95.6 \pm 0.5
	<i>leak</i>	51.5	99.1	85.0 \pm 11.6	74.6 \pm 11.1
MNLI	<i>acc</i>	49.2	91.9	—	86.9 \pm 0.5
	<i>leak</i>	34.2	91.5	—	37.4 \pm 0.7

Table 2: Accuracy and privacy metrics.
Flan-T5-Large.

Dataset		Without LoRAs	Regular FT	DC	P ³ EFT
SST2	<i>acc</i>	92.8	96.1	95.0 \pm 0.1	96.1 \pm 0.1
	<i>leak</i>	55.8	98.3	68.1 \pm 5.0	74.1 \pm 3.0
QNLI	<i>acc</i>	83.2	95.3	95.2 \pm 0.1	94.7 \pm 0.0
	<i>leak</i>	58.7	98.9	67.0 \pm 1.2	63.0 \pm 0.8
MNLI	<i>acc</i>	73.9	90.5	89.8 \pm 0.1	90.1 \pm 0.1
	<i>leak</i>	34.6	85.9	45.6 \pm 0.8	40.0 \pm 1.1

328 The results for DeBERTa are presented in Table 1. To improve reproducibility, we reuse the hyperpa-
 329 rameters from original paper, with the exception of the LoRA dropout value. We disable dropout
 330 because it interferes with the mixing weights (5). In preliminary experiments, we observed that with
 331 dropout enabled, both our algorithm and DC begin to perform significantly worse.

332 We use $n = 2$ adapter sets for P³EFT for all datasets and adhered to the same approach for the
 333 other models as well. Overall, P³FT achieves nearly the same accuracy as traditional (non-private)
 334 fine-tuning, outperforming the DC-based algorithm in terms of accuracy given the same privacy level.
 335 On the MNLI dataset, we could not find the hyperparameters for DC that ensure stable training while
 336 maintaining privacy. Meanwhile, P³EFT maintains consistent performance on this task with a slight
 337 drop in quality.

338 Table 2 reports evaluation for the Flan-T5 base model[4]. For this model, we adapt the exact same
 339 hyperparameters as in the previous evaluation with DeBERTa-XXLarge. Compared to DeBERTa,
 340 these results are more closely matched. Both both our algorithm and DC consistently solve all three
 341 tasks, but P³EFT slightly outperforms DC in terms of privacy.

Table 3: Accuracy and privacy metrics for LLaMA-2 7B.

Dataset		Without LoRAs	Regular FT	DC	P ³ EFT
SST2	<i>acc</i>	94.6	97.4	97.1 \pm 0.1	95.8 \pm 0.1
	<i>leak</i>	59.1	99.3	83.6 \pm 10.6	68.9 \pm 2.6
QNLI	<i>acc</i>	77.0	95.0	95.2 \pm 0.1	94.7 \pm 0.2
	<i>leak</i>	53.3	85.5	66.6 \pm 4.1	62.9 \pm 0.8

342 To evaluate how our algorithm scales to larger models, we also fine-tune Llama-2 7B [35] on
 343 SST2 [32] and QNLI [39] datasets. For these evaluations, we use LoRA hyperparameters that Hu
 344 et al. [15] used when fine-tuning GPT-3, with several changes inspired by Dettmers et al. [5]. Namely,
 345 we use the NF4 weight format, apply LoRA to both attention and MLP layers with rank 16. We
 346 fine-tune both tasks with maximum context length of 512 and weight decay 0.01. Table 3 summarizes
 347 our results: for QNLI, P³EFT achieves somewhat better privacy-accuracy trade-off. On SST2, P³EFT
 348 shows similarly favorable trade-offs while DC struggles to preserve privacy.

349 5 Conclusion and Discussion

350 In this work, we analyze privacy-preserving fine-tuning of large neural networks in the context of
 351 parameter-efficient fine-tuning and the two-party split learning setting. We show that while standard
 352 fine-tuning suffers from label leakage even in the parameter-efficient case, it is possible to leverage
 353 the efficiency of PEFT to alter the procedure without any significant performance drawbacks. We test
 354 the resulting method, named P³EFT, on a range of pretrained language models and multiple datasets,
 355 showing that it is competitive with a strong baseline in terms of label privacy while having higher
 356 task performance.

357 In future work, it is natural to explore how this approach can be extended to establish holistic privacy
 358 in both labels and inputs. This problem can be approached from two directions: either adapt the
 359 ideas of P³EFT for input privacy, or combine it with an existing work like [22]. Another important
 360 direction for future research is exploring the privacy of the long-term client-provider interaction. In a
 361 typical real-world use case of API fine-tuning, a client performs multiple training runs on overlapping
 362 data and hyperparameters. This could open additional attacks vectors that combine information from
 363 multiple training runs.

364 **References**

- 365 [1] Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan,
366 Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for
367 privacy-preserving machine learning. In *proceedings of the 2017 ACM SIGSAC Conference on*
368 *Computer and Communications Security*, pages 1175–1191, 2017.
- 369 [2] Alexander Borzunov, Dmitry Baranchuk, Tim Dettmers, Max Ryabinin, Younes Belkada, Artem
370 Chumachenko, Pavel Samygin, and Colin Raffel. Petals: Collaborative inference and fine-tuning
371 of large models. *arXiv preprint arXiv:2209.01188*, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2209.01188)
372 [2209.01188](https://arxiv.org/abs/2209.01188).
- 373 [3] Tianqi Chen and Carlos Guestrin. XGBoost: A scalable tree boosting system. In *Proceedings of*
374 *the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*,
375 KDD '16, pages 785–794, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4232-2. doi:
376 [10.1145/2939672.2939785](https://doi.org/10.1145/2939672.2939785). URL <http://doi.acm.org/10.1145/2939672.2939785>.
- 377 [4] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li,
378 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu,
379 Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav
380 Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H.
381 Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling
382 instruction-finetuned language models, 2022. URL <https://arxiv.org/abs/2210.11416>.
- 383 [5] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient
384 finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- 385 [6] Dreambooth API. Dreambooth API – Easily finetune Stable Diffusion and generate customised
386 AI images — dreamboothapi.ai. <https://dreamboothapi.ai/>, 2023. [Accessed 28-09-
387 2023].
- 388 [7] Haonan Duan, Adam Dziedziec, Nicolas Papernot, and Franziska Boenisch. Flocks of stochastic
389 parrots: Differentially private prompt learning for large language models. *arXiv preprint*
390 *arXiv:2305.15594*, 2023.
- 391 [8] Cynthia Dwork. Differential privacy. In *International colloquium on automata, languages, and*
392 *programming*, pages 1–12. Springer, 2006.
- 393 [9] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation.
394 In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference*
395 *on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1180–
396 1189, Lille, France, 07–09 Jul 2015. PMLR. URL [https://proceedings.mlr.press/v37/](https://proceedings.mlr.press/v37/ganin15.html)
397 [ganin15.html](https://proceedings.mlr.press/v37/ganin15.html).
- 398 [10] Otkrist Gupta and Ramesh Raskar. Distributed learning of deep neural network over multiple
399 agents. *Journal of Network and Computer Applications*, 116:1–8, 2018. ISSN 1084-8045.
400 doi: <https://doi.org/10.1016/j.jnca.2018.05.003>. URL [https://www.sciencedirect.com/](https://www.sciencedirect.com/science/article/pii/S1084804518301590)
401 [science/article/pii/S1084804518301590](https://www.sciencedirect.com/science/article/pii/S1084804518301590).
- 402 [11] Karen Hambardzumyan, Hrant Khachatrian, and Jonathan May. WARP: Word-level Ad-
403 versarial ReProgramming. In *Proceedings of the 59th Annual Meeting of the Association*
404 *for Computational Linguistics and the 11th International Joint Conference on Natural*
405 *Language Processing (Volume 1: Long Papers)*, pages 4921–4933, Online, August 2021.
406 Association for Computational Linguistics. doi: [10.18653/v1/2021.acl-long.381](https://doi.org/10.18653/v1/2021.acl-long.381). URL
407 <https://aclanthology.org/2021.acl-long.381>.
- 408 [12] Stephen Hardy, Wilko Henecka, Hamish Ivey-Law, Richard Nock, Giorgio Patrini, Guillaume
409 Smith, and Brian Thorne. Private federated learning on vertically partitioned data via entity
410 resolution and additively homomorphic encryption, 2017.
- 411 [13] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. DeBERTa: Decoding-enhanced
412 bert with disentangled attention. In *International Conference on Learning Representations*,
413 2021. URL <https://openreview.net/forum?id=XPZiaotutsD>.

- 414 [14] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe,
415 Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning
416 for NLP. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th*
417 *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning*
418 *Research*, pages 2790–2799. PMLR, 09–15 Jun 2019. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v97/houlsby19a.html)
419 [press/v97/houlsby19a.html](https://proceedings.mlr.press/v97/houlsby19a.html).
- 420 [15] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang,
421 Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In
422 *International Conference on Learning Representations*, 2022. URL [https://openreview.](https://openreview.net/forum?id=nZeVKeeFYf9)
423 [net/forum?id=nZeVKeeFYf9](https://openreview.net/forum?id=nZeVKeeFYf9).
- 424 [16] Hugging Face. AutoTrain — huggingface.co. <https://huggingface.co/autotrain>, 2023.
425 [Accessed 28-09-2023].
- 426 [17] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
427 *arXiv:1412.6980*, 2014.
- 428 [18] Oscar Li, Jiankai Sun, Xin Yang, Weihao Gao, Hongyi Zhang, Junyuan Xie, Virginia Smith,
429 and Chong Wang. Label leakage and protection in two-party split learning. In *International*
430 *Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=c0tBRgsf2f0)
431 [id=c0tBRgsf2f0](https://openreview.net/forum?id=c0tBRgsf2f0).
- 432 [19] Qinbin Li, Zeyi Wen, Zhaomin Wu, Sixu Hu, Naibo Wang, Yuan Li, Xu Liu, and Bingsheng He.
433 A survey on federated learning systems: Vision, hype and reality for data privacy and protection.
434 *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- 435 [20] Shen Li, Pritam Damania, Luca Wehrstedt, and Rohan Varma. PyTorch RPC: Distributed Deep
436 Learning Built on Tensor-Optimized Remote Procedure Calls. In *Proceedings of Machine*
437 *Learning and Systems 5 (MLSys)*, 2023.
- 438 [21] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In
439 *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the*
440 *11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*,
441 pages 4582–4597, Online, August 2021. Association for Computational Linguistics. doi:
442 [10.18653/v1/2021.acl-long.353](https://doi.org/10.18653/v1/2021.acl-long.353). URL <https://aclanthology.org/2021.acl-long.353>.
- 443 [22] Yansong Li, Zhixing Tan, and Yang Liu. Privacy-preserving prompt tuning for large language
444 model services. *ArXiv*, abs/2305.06212, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:258588141)
445 [CorpusID:258588141](https://api.semanticscholar.org/CorpusID:258588141).
- 446 [23] Xiao-Yang Liu, Rongyi Zhu, Daochen Zha, Jiechao Gao, Shan Zhong, and Meikang Qiu.
447 Differentially private low-rank adaptation of large language model using federated learning.
448 *arXiv preprint arXiv:2312.17493*, 2023.
- 449 [24] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Ar-
450 cas. Communication-Efficient Learning of Deep Networks from Decentralized Data. In Aarti
451 Singh and Jerry Zhu, editors, *Proceedings of the 20th International Conference on Artificial*
452 *Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*,
453 pages 1273–1282. PMLR, 20–22 Apr 2017. URL [https://proceedings.mlr.press/v54/](https://proceedings.mlr.press/v54/mcmahan17a.html)
454 [mcmahan17a.html](https://proceedings.mlr.press/v54/mcmahan17a.html).
- 455 [25] Nvidia. Nvidia confidential computing. [https://www.nvidia.com/en-us/data-center/](https://www.nvidia.com/en-us/data-center/solutions/confidential-computing)
456 [solutions/confidential-computing](https://www.nvidia.com/en-us/data-center/solutions/confidential-computing), 2023. [Accessed 28-09-2023].
- 457 [26] OctoAI. Fine-tuning Stable Diffusion — docs.octoai.cloud. [https://docs.octoai.cloud/](https://docs.octoai.cloud/docs/fine-tuning-stable-diffusion)
458 [docs/fine-tuning-stable-diffusion](https://docs.octoai.cloud/docs/fine-tuning-stable-diffusion), 2023. [Accessed 28-09-2023].
- 459 [27] OpenAI. OpenAI Platform — platform.openai.com. [https://platform.openai.com/](https://platform.openai.com/docs/guides/fine-tuning)
460 [docs/guides/fine-tuning](https://platform.openai.com/docs/guides/fine-tuning), 2023. [Accessed 28-09-2023].

- 461 [28] Dario Pasquini, Giuseppe Ateniese, and Massimo Bernaschi. Unleashing the tiger: Inference
462 attacks on split learning. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and*
463 *Communications Security, CCS '21*, page 2113–2129, New York, NY, USA, 2021. Association
464 for Computing Machinery. ISBN 9781450384544. doi: 10.1145/3460120.3485259. URL
465 <https://doi.org/10.1145/3460120.3485259>.
- 466 [29] Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych.
467 Adapterfusion: Non-destructive task composition for transfer learning, 2021.
- 468 [30] Yuma Rao, Jacob Steeves, Ala Shaabana, Daniel Attevelt, and Matthew McAteer. Bittensor: A
469 peer-to-peer intelligence market, 2021.
- 470 [31] Weiyan Shi, Ryan Shea, Si Chen, Chiyuan Zhang, Ruoxi Jia, and Zhou Yu. Just fine-tune twice:
471 Selective differential privacy for large language models. *arXiv preprint arXiv:2204.07667*,
472 2022.
- 473 [32] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew
474 Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a
475 sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural*
476 *Language Processing*, pages 1631–1642, Seattle, Washington, USA, October 2013. Association
477 for Computational Linguistics. URL <https://www.aclweb.org/anthology/D13-1170>.
- 478 [33] Jiankai Sun, Xin Yang, Yuanshun Yao, and Chong Wang. Label leakage and protection from
479 forward embedding in vertical federated learning. *arXiv preprint arXiv:2203.01451*, 2022.
- 480 [34] Chuanqi Tan, Fuchun Sun, Tao Kong, Wenchang Zhang, Chao Yang, and Chunfang Liu. A
481 survey on deep transfer learning. In Věra Kůrková, Yannis Manolopoulos, Barbara Hammer,
482 Lazaros Iliadis, and Ilias Maglogiannis, editors, *Artificial Neural Networks and Machine*
483 *Learning – ICANN 2018*, pages 270–279, Cham, 2018. Springer International Publishing. ISBN
484 978-3-030-01424-7.
- 485 [35] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
486 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open
487 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 488 [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N
489 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon,
490 U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, ed-
491 itors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,
492 Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/file/](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
493 [3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 494 [37] Praneeth Vepakomma, Otkrist Gupta, Tristan Swedish, and Ramesh Raskar. Split learning for
495 health: Distributed deep learning without sharing raw patient data, 2018.
- 496 [38] Praneeth Vepakomma, Otkrist Gupta, Abhimanyu Dubey, and Ramesh Raskar. Reducing
497 leakage in distributed deep learning for sensitive health data. 05 2019.
- 498 [39] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.
499 Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv*
500 *preprint arXiv:1804.07461*, 2018.
- 501 [40] Yiming Wang, Yu Lin, Xiaodong Zeng, and Guannan Zhang. Privatelora for efficient privacy
502 preserving llm. *arXiv preprint arXiv:2311.14030*, 2023.
- 503 [41] Adina Williams, Nikita Nangia, and Samuel R Bowman. A broad-coverage challenge corpus
504 for sentence understanding through inference. *arXiv preprint arXiv:1704.05426*, 2017.
- 505 [42] Guangxuan Xiao, Ji Lin, and Song Han. Offsite-tuning: Transfer learning without full model.
506 *arXiv preprint arXiv:2302.04870*, 2023.
- 507 [43] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept
508 and applications. *ACM Trans. Intell. Syst. Technol.*, 10(2), jan 2019. ISSN 2157-6904. doi:
509 10.1145/3298981. URL <https://doi.org/10.1145/3298981>.

- 510 [44] Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath,
 511 Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, Sergey Yekhanin, and
 512 Huishuai Zhang. Differentially private fine-tuning of language models. In *International
 513 Conference on Learning Representations, 2022*. URL [https://openreview.net/forum/
 514 id=Q42f0dfjEC0](https://openreview.net/forum?id=Q42f0dfjEC0).
- 515 [45] Zhuo Zhang, Yuanhang Yang, Yong Dai, Qifan Wang, Yue Yu, Lizhen Qu, and Zenglin
 516 Xu. FedPETuning: When federated learning meets the parameter-efficient tuning methods of
 517 pre-trained language models. In *Findings of the Association for Computational Linguistics:
 518 ACL 2023*, pages 9963–9977, Toronto, Canada, July 2023. Association for Computational
 519 Linguistics. doi: 10.18653/v1/2023.findings-acl.632. URL [https://aclanthology.org/
 520 2023.findings-acl.632](https://aclanthology.org/2023.findings-acl.632).
- 521 [46] Haodong Zhao, Wei Du, Fangqi Li, Peixuan Li, and Gongshen Liu. Fedprompt: Communication-
 522 efficient and privacy preserving prompt tuning in federated learning, 2023.

523 A Hyperparameters search

524 In P³EFT and Distance Correlation methods resulting loss L function can be viewed in the form

$$L = L_m + \alpha \cdot L_r,$$

525 where L_m - main task loss, L_r - regularizer and α is a coefficient that controls the tradeoff between
 526 these two losses. The selection of this coefficient affects the final performance of the model. Therefore,
 527 to find the best configurations for both methods, we iterated through this hyperparameter using a grid
 528 search.

529 We started with $\alpha = 1$ and then altered it with a multiplicative step of $10^{\frac{1}{2}}$. Values were discarded if
 530 the quality did not exceed that achieved by solely training the classifier without LoRA. This criterion
 531 was adopted because such outcomes would suggest the method’s inability to outperform training
 532 scenarios in which the server does not engage with the labels whatsoever. Additionally, we excluded
 533 values that led to unstable training. By this, we mean instances where, although the model initially
 534 trained on the primary task, at some point, the regularizer began contributing significantly more,
 535 and the utility value dropped to the starting value. We observed this issue for the DC method with
 536 DeBERTa on the MNLI. From the remaining values, we aimed to choose the one that offered the
 537 lowest privacy leakage. The final hyperparameter values for P³EFT can be found in the Table 4 and
 538 for DC in the Table 5.

Table 4: Regularization parameter α for the P³EFT method. The values in the table represent powers of the $10^{\frac{1}{2}}$.

	SST2	QNLI	MNLI
DeBERTa XXLlarge	1	1	1
Flan-T5-Large	-1	1	1
LLaMA-2 7B	0	0	—

Table 5: Regularization parameter α for the DC method. The values in the table represent powers of the $10^{\frac{1}{2}}$.

	SST2	QNLI	MNLI
DeBERTa XXLlarge	0	-1	—
Flan-T5-Large	2	-1	0
LLaMA-2 7B	-1	-1	—

539 B Formal algorithm definition

540 Below, we define the full P³EFT algorithm. In Algorithm 2, main_loss is the task-specific objective
 541 e.g. cross-entropy; reg_loss is the adversarial regularizer described in Section 3.4. We denote
 542 client-side model "head" as $f(h, \psi^t)$, where ψ represent trainable head parameters. Finally, opt_step
 543 function performs a single gradient descent step with a task-specific optimizer, typically Adam [17].

Algorithm 2 P³EFT - full training algorithm

```

1: Input: dataset  $D = \{X, Y\}$ ,  $n > 1$  number of adapters,  $\alpha \geq 0$  - regularizing weight,  $m > 1$ 
   number of obfuscated gradients
2: Initialize head  $\psi^0$ , mixing weights  $W_i$  and adapters  $\theta_i^0, i = \overline{1, n}$ 
3: for  $t = 0, 1, \dots, T - 1$  do
4:   Sample batch  $\{x^t, y^t\}$ 
5:   for  $i = 1, \dots, n$  do
6:      $h_i^t = h(x^t, \theta_i^t)$  ▷ by server
7:      $l_i = \text{reg\_loss}(h_i^t, y^t)$  ▷ by client
8:   end for
9:    $h' = \sum_{i=1}^n W_i \odot h_i^t$ 
10:   $l = \text{main\_loss}(f(h', \psi^t), y^t)$ 
11:   $L = l + \alpha \cdot \sum_{i=1}^n l_i$ 
12:  for  $i = 1, \dots, n$  do ▷ Client performs partial backprop
13:     $g_h = \partial L / \partial h_i^t$ 
14:     $g_i^t = \text{private\_backprop}(x, \theta_i^t, g_h, m)$ 
15:     $\theta_i^{t+1} = \text{opt\_step}(\theta_i^t, g_i^t, t)$ 
16:  end for
17:   $\psi^{t+1} = \text{opt\_step}(\psi^t, \partial l / \partial \psi^t, t)$ 
18: end for
19: Return:  $\psi^T, \theta_1^T, \dots, \theta_M^T$ 

```

544 C Informal description of LoRA fine-tuning

545 For convenience, we provide a brief summary of fine-tuning with LoRA [15]. This PEFT method
 546 was originally designed for fine-tuning large pre-trained language models on downstream NLP tasks.
 547 These language models are typically based on the Transformer architecture [36], where most trainable
 548 parameters are allocated to linear layers in multi-head self-attention and feedforward blocks.

549 In the first stage of LoRA fine-tuning, user augments the model with adapters. To do so, a user goes
 550 over linear layers in transformer blocks and adds two trainable matrices, A and B that affect this
 551 layer's forward pass. Let $W_i \times x + b_i$ be the original layer with n inputs and m hidden units. Here,
 552 $W_i \in \mathcal{R}^{m \times n}$ is a pre-trained weight matrix, $b_i \in \mathcal{R}^m$ is a pre-trained intercept vector and $x \in \mathcal{R}^n$
 553 represents a vector of inputs to this particular layer. During the forward pass, a layer with LoRA
 554 adapter computes $W_i \times x + b_i + B_i \times A_i \times x$, or equivalently, $(W_i + B \times A) \times x + b_i$. Here, A_i
 555 and B_i are two newly added matrices that constitute a LoRA adapter.

556 The adapter matrices $A \in \mathcal{R}^{r \times n}$ and $B \in \mathcal{R}^{m \times r}$ have a very small intermediate dimension r . For
 557 instance, when training GPT-3 with LoRA adapters, authors use $1 \leq r \leq 64$, whereas the main
 558 weight dimensions are $m = n = 12288$. The first matrix A is initialized with small random normal
 559 values, and the second matrix B is initialized at zeros. That way, initial A and B do not affect the
 560 model predictions.

561 Once all adapters are initialized, the user trains all A_i and B_i matrices of the model, while keeping
 562 the rest of the weights frozen. This way, only a small fraction (less than 1%) of model weights are
 563 updated. Once the training is over, the learned adapters A_i and B_i can be merged into the main
 564 weights ($W_i := W_i + A_i \times B_i$) or used separately.

565 LoRA adapters are designed with two objectives in mind: i) to allow fine-tuning models in limited
 566 GPU memory and ii) to allow inferencing many fine-tuned models using one inference server. When
 567 fine-tuning, LoRA achieves small memory footprint due to the fact that user does not need to compute

568 gradients (or optimizer statistics) for billions of main model parameters. During inference, a server
569 can keep a library of several adapters for different tasks and swap between them on demand.

570 **D Informal description of LoRA fine-tuning**

571 We used NVIDIA A100 GPUs for all the experiments. Experiments with DeBERTA [13] and Flan-T5
572 [4] on SST2 [32] were conducted on the single GPU, while experiments on MNLI [41] and QNLI
573 require 4 A100. LLaMA-2 [35] experiments were carried out on the node of 8 A100.

574 All the experiments last 12-24 hours. However, it is possible to speed up some of them using more
575 GPUs, as well as conduct them on a smaller number of GPUs using technics to save GPU memory
576 (see parameters in our code).

577 **NeurIPS Paper Checklist**

578 **1. Claims**

579 Question: Do the main claims made in the abstract and introduction accurately reflect the
580 paper's contributions and scope?

581 Answer: [Yes]

582 Justification: See our main contributions in 1 and Section 5.

583 Guidelines:

- 584 • The answer NA means that the abstract and introduction do not include the claims
585 made in the paper.
- 586 • The abstract and/or introduction should clearly state the claims made, including the
587 contributions made in the paper and important assumptions and limitations. A No or
588 NA answer to this question will not be perceived well by the reviewers.
- 589 • The claims made should match theoretical and experimental results, and reflect how
590 much the results can be expected to generalize to other settings.
- 591 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
592 are not attained by the paper.

593 **2. Limitations**

594 Question: Does the paper discuss the limitations of the work performed by the authors?

595 Answer: [Yes]

596 Justification: We discuss potential limitations of the method in Section 1, Section 3 and
597 Section 5.

598 Guidelines:

- 599 • The answer NA means that the paper has no limitation while the answer No means that
600 the paper has limitations, but those are not discussed in the paper.
- 601 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 602 • The paper should point out any strong assumptions and how robust the results are to
603 violations of these assumptions (e.g., independence assumptions, noiseless settings,
604 model well-specification, asymptotic approximations only holding locally). The authors
605 should reflect on how these assumptions might be violated in practice and what the
606 implications would be.
- 607 • The authors should reflect on the scope of the claims made, e.g., if the approach was
608 only tested on a few datasets or with a few runs. In general, empirical results often
609 depend on implicit assumptions, which should be articulated.
- 610 • The authors should reflect on the factors that influence the performance of the approach.
611 For example, a facial recognition algorithm may perform poorly when image resolution
612 is low or images are taken in low lighting. Or a speech-to-text system might not be
613 used reliably to provide closed captions for online lectures because it fails to handle
614 technical jargon.
- 615 • The authors should discuss the computational efficiency of the proposed algorithms
616 and how they scale with dataset size.
- 617 • If applicable, the authors should discuss possible limitations of their approach to
618 address problems of privacy and fairness.
- 619 • While the authors might fear that complete honesty about limitations might be used by
620 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
621 limitations that aren't acknowledged in the paper. The authors should use their best
622 judgment and recognize that individual actions in favor of transparency play an impor-
623 tant role in developing norms that preserve the integrity of the community. Reviewers
624 will be specifically instructed to not penalize honesty concerning limitations.

625 **3. Theory Assumptions and Proofs**

626 Question: For each theoretical result, does the paper provide the full set of assumptions and
627 a complete (and correct) proof?

628 Answer: [NA]

629 Justification: We include no proofs.

630 Guidelines:

- 631 • The answer NA means that the paper does not include theoretical results.
- 632 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 633 referenced.
- 634 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 635 • The proofs can either appear in the main paper or the supplemental material, but if
- 636 they appear in the supplemental material, the authors are encouraged to provide a short
- 637 proof sketch to provide intuition.
- 638 • Inversely, any informal proof provided in the core of the paper should be complemented
- 639 by formal proofs provided in appendix or supplemental material.
- 640 • Theorems and Lemmas that the proof relies upon should be properly referenced.

641 4. Experimental Result Reproducibility

642 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

643 perimental results of the paper to the extent that it affects the main claims and/or conclusions

644 of the paper (regardless of whether the code and data are provided or not)?

645 Answer: [Yes]

646 Justification: We describe all the technical details in the Section 4 and in the README of

647 our attached code.

648 Guidelines:

- 649 • The answer NA means that the paper does not include experiments.
- 650 • If the paper includes experiments, a No answer to this question will not be perceived
- 651 well by the reviewers: Making the paper reproducible is important, regardless of
- 652 whether the code and data are provided or not.
- 653 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 654 to make their results reproducible or verifiable.
- 655 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 656 For example, if the contribution is a novel architecture, describing the architecture fully
- 657 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 658 be necessary to either make it possible for others to replicate the model with the same
- 659 dataset, or provide access to the model. In general, releasing code and data is often
- 660 one good way to accomplish this, but reproducibility can also be provided via detailed
- 661 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 662 of a large language model), releasing of a model checkpoint, or other means that are
- 663 appropriate to the research performed.
- 664 • While NeurIPS does not require releasing code, the conference does require all submis-
- 665 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 666 nature of the contribution. For example
- 667 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 668 to reproduce that algorithm.
- 669 (b) If the contribution is primarily a new model architecture, the paper should describe
- 670 the architecture clearly and fully.
- 671 (c) If the contribution is a new model (e.g., a large language model), then there should
- 672 either be a way to access this model for reproducing the results or a way to reproduce
- 673 the model (e.g., with an open-source dataset or instructions for how to construct
- 674 the dataset).
- 675 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 676 authors are welcome to describe the particular way they provide for reproducibility.
- 677 In the case of closed-source models, it may be that access to the model is limited in
- 678 some way (e.g., to registered users), but it should be possible for other researchers
- 679 to have some path to reproducing or verifying the results.

680 5. Open access to data and code

681 Question: Does the paper provide open access to the data and code, with sufficient instruc-

682 tions to faithfully reproduce the main experimental results, as described in supplemental

683 material?

684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735

Answer: [Yes]

Justification: We provide attached code.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See Section 4, Appendix A.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We report error bars in main results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 736
- 737
- 738
- 739
- 740
- 741
- 742
- 743
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

744 **8. Experiments Compute Resources**

745 Question: For each experiment, does the paper provide sufficient information on the com-
746 puter resources (type of compute workers, memory, time of execution) needed to reproduce
747 the experiments?

748 Answer: [Yes]

749 Justification: See Appendix D.

750 Guidelines:

- 751
- 752
- 753
- 754
- 755
- 756
- 757
- 758
- The answer NA means that the paper does not include experiments.
 - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
 - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

759 **9. Code Of Ethics**

760 Question: Does the research conducted in the paper conform, in every respect, with the
761 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

762 Answer: [Yes]

763 Justification: We confirm the NeurIPS Code of Ethics

764 Guidelines:

- 765
- 766
- 767
- 768
- 769
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

770 **10. Broader Impacts**

771 Question: Does the paper discuss both potential positive societal impacts and negative
772 societal impacts of the work performed?

773 Answer: [Yes]

774 Justification: We describe a potential social impact in Introduction.

775 Guidelines:

- 776
- 777
- 778
- 779
- 780
- 781
- 782
- 783
- 784
- 785
- 786
- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to

787 generate deepfakes for disinformation. On the other hand, it is not needed to point out
788 that a generic algorithm for optimizing neural networks could enable people to train
789 models that generate Deepfakes faster.

- 790 • The authors should consider possible harms that could arise when the technology is
791 being used as intended and functioning correctly, harms that could arise when the
792 technology is being used as intended but gives incorrect results, and harms following
793 from (intentional or unintentional) misuse of the technology.
- 794 • If there are negative societal impacts, the authors could also discuss possible mitigation
795 strategies (e.g., gated release of models, providing defenses in addition to attacks,
796 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
797 feedback over time, improving the efficiency and accessibility of ML).

798 11. Safeguards

799 Question: Does the paper describe safeguards that have been put in place for responsible
800 release of data or models that have a high risk for misuse (e.g., pretrained language models,
801 image generators, or scraped datasets)?

802 Answer: [No]

803 Justification: We do not describe the safeguards

804 Guidelines:

- 805 • The answer NA means that the paper poses no such risks.
- 806 • Released models that have a high risk for misuse or dual-use should be released with
807 necessary safeguards to allow for controlled use of the model, for example by requiring
808 that users adhere to usage guidelines or restrictions to access the model or implementing
809 safety filters.
- 810 • Datasets that have been scraped from the Internet could pose safety risks. The authors
811 should describe how they avoided releasing unsafe images.
- 812 • We recognize that providing effective safeguards is challenging, and many papers do
813 not require this, but we encourage authors to take this into account and make a best
814 faith effort.

815 12. Licenses for existing assets

816 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
817 the paper, properly credited and are the license and terms of use explicitly mentioned and
818 properly respected?

819 Answer: [Yes]

820 Justification: We use open datasets from GLUE benchmark and open-sourced models and
821 cite them in our work.

822 Guidelines:

- 823 • The answer NA means that the paper does not use existing assets.
- 824 • The authors should cite the original paper that produced the code package or dataset.
- 825 • The authors should state which version of the asset is used and, if possible, include a
826 URL.
- 827 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 828 • For scraped data from a particular source (e.g., website), the copyright and terms of
829 service of that source should be provided.
- 830 • If assets are released, the license, copyright information, and terms of use in the
831 package should be provided. For popular datasets, paperswithcode.com/datasets
832 has curated licenses for some datasets. Their licensing guide can help determine the
833 license of a dataset.
- 834 • For existing datasets that are re-packaged, both the original license and the license of
835 the derived asset (if it has changed) should be provided.
- 836 • If this information is not available online, the authors are encouraged to reach out to
837 the asset's creators.

838 13. New Assets

839 Question: Are new assets introduced in the paper well documented and is the documentation
840 provided alongside the assets?

841 Answer: [NA]

842 Justification: We do not release any of the assets

843 Guidelines:

- 844 • The answer NA means that the paper does not release new assets.
- 845 • Researchers should communicate the details of the dataset/code/model as part of their
846 submissions via structured templates. This includes details about training, license,
847 limitations, etc.
- 848 • The paper should discuss whether and how consent was obtained from people whose
849 asset is used.
- 850 • At submission time, remember to anonymize your assets (if applicable). You can either
851 create an anonymized URL or include an anonymized zip file.

852 14. Crowdsourcing and Research with Human Subjects

853 Question: For crowdsourcing experiments and research with human subjects, does the paper
854 include the full text of instructions given to participants and screenshots, if applicable, as
855 well as details about compensation (if any)?

856 Answer: [NA]

857 Justification: We do not involve crowdsourcing.

858 Guidelines:

- 859 • The answer NA means that the paper does not involve crowdsourcing nor research with
860 human subjects.
- 861 • Including this information in the supplemental material is fine, but if the main contribu-
862 tion of the paper involves human subjects, then as much detail as possible should be
863 included in the main paper.
- 864 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
865 or other labor should be paid at least the minimum wage in the country of the data
866 collector.

867 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 868 Subjects

869 Question: Does the paper describe potential risks incurred by study participants, whether
870 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
871 approvals (or an equivalent approval/review based on the requirements of your country or
872 institution) were obtained?

873 Answer: [NA]

874 Justification: Our paper does not involve research with human subjects

875 Guidelines:

- 876 • The answer NA means that the paper does not involve crowdsourcing nor research with
877 human subjects.
- 878 • Depending on the country in which research is conducted, IRB approval (or equivalent)
879 may be required for any human subjects research. If you obtained IRB approval, you
880 should clearly state this in the paper.
- 881 • We recognize that the procedures for this may vary significantly between institutions
882 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
883 guidelines for their institution.
- 884 • For initial submissions, do not include any information that would break anonymity (if
885 applicable), such as the institution conducting the review.