

---

# Voting-based Approaches for Differentially Private Federated Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Differentially Private Federated Learning (DPFL) is an emerging field with many  
2 applications. Gradient averaging based DPFL methods require costly commu-  
3 nication rounds and hardly work with large-capacity models, due to the explicit  
4 dimension dependence in its added noise. In this paper, inspired by the non-  
5 federated knowledge transfer privacy learning methods, we design two DPFL  
6 algorithms (*AE-DPFL* and *kNN-DPFL*) that provide provable DP guarantees for  
7 both instance-level and agent-level privacy regimes. By voting among the data  
8 *labels* returned from each local model, instead of averaging the gradients, our  
9 algorithms avoid the dimension dependence and significantly reduces the commu-  
10 nication cost. Theoretically, by applying secure multi-party computation, we could  
11 exponentially amplify the (data-dependent) privacy guarantees when the margin  
12 of the voting scores are distinctive. Empirical evaluation on both instance and  
13 agent level DP is conducted across five datasets. When aligning privacy cost the  
14 same, we show 2% to 12% higher accuracy compared to DP-FedAvg, or aligning  
15 accuracy the same, we show that less than 65% privacy cost is achieved.

## 16 1 Introduction

17 Federated learning (FL) [McMahan et al., 2017, Bonawitz et al., 2017b, Mohassel and Zhang, 2017,  
18 Smith et al., 2017] is an emerging paradigm of distributed machine learning with a wide range of  
19 applications [Kairouz et al., 2019]. FL allows distributed agents to collaboratively train a centralized  
20 machine learning model without sharing each of their local data, thereby sidestepping the ethical and  
21 legal concerns that arise in collecting private user data for the purpose of building machine-learning  
22 based products and services.

23 The workflow of FL is often enhanced by secure multi-party computation [Bonawitz et al., 2017b]  
24 (MPC) so as to handle various threat models in the communication protocols, which provably ensures  
25 that agents can receive the output of the computation (e.g., the sum of the gradients) but nothing in  
26 between (e.g., other agents' gradients).

27 However, MPC alone does not protect the agents or their users from inference attacks that use only  
28 the output, or combine the output with auxiliary information. Extensive studies demonstrate that  
29 these attacks may lead to a blatant reconstruction of proprietary datasets [Dinur and Nissim, 2003],  
30 high-confidence identification of individuals (a legal liability for the participating agents) [Shokri  
31 et al., 2017], or even completion of social security numbers [Carlini et al., 2019]. Motivated by  
32 these challenges, there have been a number of recent efforts [Truex et al., 2019, Geyer et al., 2017,  
33 McMahan et al., 2018] in developing federated learning methods with differential privacy (DP)  
34 [Dwork et al., 2006], which is a well-established definition of privacy that provably prevents such  
35 attacks.

36 Existing methods in differentially private federated learning (DPFL), e.g., DP-FedAvg [Geyer et al.,  
37 2017, McMahan et al., 2018] and the recent state-of-the-art DP-FedSGD [Truex et al., 2019], are

38 predominantly noisy gradient based methods, which build upon the NoisySGD method, a classical  
39 algorithm in (non-federated) DP learning [Song et al., 2013, Bassily et al., 2014, Abadi et al.,  
40 2016]. They work by iteratively aggregating (multi-)gradient updates from individual agents using a  
41 differentially private mechanism. A notable limitation for this approach is that they require clipping  
42 the  $\ell_2$  magnitude of gradients to a threshold  $S$  and adding noise proportional to  $S$  to *every coordinate*  
43 of the high dimensional parameters from the shared global model. The clipping and perturbation steps  
44 introduce either large bias (when  $S$  is small) or large variance (when  $S$  is large), which interferes with  
45 convergence of SGD, which makes scaling to large-capacity models difficult. In Sec. A, we concretely  
46 demonstrate these limitations with examples and theory. Particularly, we show that FedAvg may  
47 fail to decrease the loss function using gradient clipping, and DP-FedAvg requires many outer-loop  
48 iterations (i.e., many rounds of communication to synchronize model parameters) to converge under  
49 differential privacy.

50 In this paper, we consider a fundamentally different DP learning setting known as the *Knowledge*  
51 *Transfer* model [Papernot et al., 2017] (a.k.a. the *Model-Agnostic Private learning* model [Bassily  
52 et al., 2018]). This model requires an *unlabeled* dataset to be available *in the clear*, which makes  
53 this setting slightly more restrictive. However, when such a public dataset is indeed available (it  
54 often is in federated learning with domain adaptation, see, e.g., Peterson et al. [2019], Mohri et al.  
55 [2019], Peng et al. [2019b]), it could substantially improve the privacy-utility tradeoff in DP learning  
56 [Papernot et al., 2017, 2018, Zhu et al., 2020].

57 The goal of this paper is to develop DPFL algorithms under the *knowledge transfer* model, for which  
58 we propose two algorithms (*AE-DPFL* and *kNN-DPFL*), that further develop from the *non-distributed*  
59 *Private-Aggregation-of-Teacher-Ensembles* (PATE) [Papernot et al., 2018] and *Private-kNN* [Zhu  
60 et al., 2020] to the FL setting. We discover that the distinctive characteristics of these algorithms  
61 make them *natural* and *highly desirable* for DPFL tasks. Specifically, the private aggregation is now  
62 essentially privately releasing “ballot counts” in the (one-hot) label space, instead of the parameter  
63 (gradient) space. This naturally avoids the aforementioned issues associated with high dimensionality  
64 and gradient clipping. Instead of transmitting the gradient update, transmitting the vote of the “ballot  
65 counts” tremendously reduce the communication cost. Moreover, many iterations of the model update  
66 using noise addition with SGD, leads to poor privacy guarantee, where our methods exactly avoid  
67 this and use voting on labels, thus significantly outperform the state-of-the-art DPFL methods.

68 Our contributions are summarized in four folds.

- 69 1. We construct examples to demonstrate that DP-FedAvg (a) may fail due to gradient clipping  
70 and (b) requires many rounds of communications (see Section Challenge in the appendix);  
71 while our approach naturally avoids both limitations.
- 72 2. We design two voting-based distributed algorithms that provide provable DP guarantees on  
73 both *agent-level* and *instance (of-each-agent)-level* granularity, which makes them suitable for  
74 both well-studied regimes of FL: (a) distributed learning from on-device data; (b) collaboration  
75 of a few large organizations.
- 76 3. We demonstrate “privacy-amplification by ArgMax” by a new MPC technique [Dery et al.,  
77 2019] — our proposed private voting mechanism enjoys an *exponentially stronger* (data-  
78 dependent) privacy guarantee when the “winner” wins by a large margin.
- 79 4. Extensive evaluation demonstrates that our method systematically improves the privacy-utility  
80 trade-off over DP-FedAvg and DP-FedSGD, and that our methods are more robust towards  
81 distribution-shifts across agents.

82 **A remark of our novelty.** Though *AE-DPFL* and *kNN-DPFL* are algorithmically similar to the  
83 original *PATE* [Papernot et al., 2018] and *Private-kNN* [Zhu et al., 2020], they are not the same and  
84 we facilitate them to a new problem — *federated learning*. The facilitation itself is nontrivial and  
85 requires substantial technical innovations. We highlight three challenges below.

86 To begin with, several key DP techniques that contribute to the success of PATE and Private-kNN  
87 in the standard settings are no longer applicable (e.g., privacy amplification by sampling and noisy  
88 screening). This is partially because in standard private learning, the attacker only sees the final  
89 models; but in FL, the attacker can eavesdrop in all network traffic and could be a subset of the agents  
90 themselves.

91 Moreover, PATE and Private-kNN only provide instance-level DP. We show *AE-DPFL* and *kNN-*  
92 *DPFL* also satisfy the stronger agent-level DP. *AE-DPFL*’s agent-level DP parameter is, interestingly,

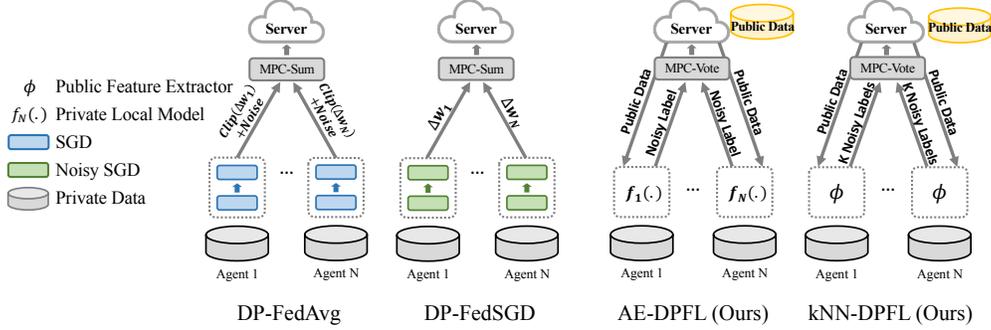


Figure 1: The structural difference of our methods to DP-FedAvg and DP-FedSGD. DP-FedAvg and AE-DPFL are for agent-level DP. DP-FedSGD and kNN-DPFL are for instance-level DP.

93 a factor of 2 better than its instance-level DP parameter. *kNN-DPFL* in addition enjoys a factor of  $k$   
 94 amplification for the instance-level DP.

95 Thirdly, a key challenge of FL is data heterogeneity of individual agents. Methods like PATE  
 96 randomly split the dataset so each teacher is identically distributed, but this assumption is violated  
 97 with heterogeneous agents. Similarly, methods like Private-kNN have also been demonstrated only  
 98 under homogeneous settings. In contrast, our proposed methods – AE-DPFL and kNN-DPFL –  
 99 exhibit robustness to data heterogeneity and domain shifts, as demonstrated in our experiments. Note  
 100 that techniques like domain adaptation may lead to further complementary benefits, but we defer its  
 101 exploration to future work, while focusing our scope here on novel techniques for DPFL.

## 102 2 Preliminary

103 Differential privacy [Dwork et al., 2006] is a quantifiable definition of privacy that provides provable  
 104 guarantees against identification of individuals in a private dataset.

105 **Definition 1. Differential Privacy:** A randomized mechanism  $\mathcal{M} : \mathcal{D} \rightarrow \mathcal{R}$  with a domain  $\mathcal{D}$  and  
 106 range  $\mathcal{R}$  satisfies  $(\epsilon, \delta)$ -differential privacy, if for any two adjacent datasets  $D, D' \in \mathcal{D}$  and for any  
 107 subset of outputs  $O \subseteq \mathcal{R}$ , it holds that  $\Pr[\mathcal{M}(D) \in O] \leq e^\epsilon \Pr[\mathcal{M}(D') \in O] + \delta$ .

108 The definition indicates that one could not distinguish between  $D$  and  $D'$  therefore protecting the  
 109 “delta” between  $D, D'$ . Depending on how *adjacency* is defined, this “delta” comes with different  
 110 semantic meaning. We consider two levels of granularity:

111 **Definition 2. Agent-level DP:** When  $D'$  is constructed by adding or removing an agent from  $D$  (with  
 112 all data points from that agent).

113 **Definition 3. Instance-level DP:** When  $D'$  is constructed by adding or removing one data point from  
 114 any of the agents.

115 The above two definitions are each important in particular situations. For example, when a smart  
 116 phone app jointly learns from its users’ text messages, it is more appropriate to protect each user as a  
 117 unit, which is agent-level DP. In another situation, when a few hospitals would like to collaborate  
 118 on a patient study through federated learning, obfuscating the entire dataset from one hospital is  
 119 meaningless, which makes instance-level DP better-suited to protect an individual patient from being  
 120 identified.

121 **DPFL Baselines:** DP-FedAvg [Geyer et al., 2017, McMahan et al., 2018] (Figure 1), a representative  
 122 DPFL algorithm, when compared to FedAvg, it enforces clipping of per-agent model gradient to a  
 123 threshold  $S$  and adds noise to the scaled gradient before it is averaged at the server, which ensures  
 124 agent-level DP. DP-FedSGD [Truex et al., 2019, Peterson et al., 2019], is one of the state-of-the-arts  
 125 that focus on instance-level DP. It performs NoisySGD [Abadi et al., 2016] for a fixed number of  
 126 iterations at each agent. The gradient updates are averaged on each communication round at the  
 127 server, as shown in Figure 1.

## 128 3 Our Approach

129 We propose two voting-base algorithms, termed aggregation ensemble DPFL “AE-DPFL” and  $k$   
 130 Nearest Neighbor DPFL “kNN-DPFL”. Each algorithm first privately labels a subset of data from the  
 131 server and then trains a global model using pseudo-labeled data.

### 132 3.1 Aggregation Ensemble - DPFL

133 In *AE-DPFL* (Algorithm 1), each agent  $i$  trains a local agent model  $f_i$  using its own private local data.  
134 The local model is never revealed to the server but only used to make predictions for unlabeled data  
135 (queries). For each query  $x_t$ , every agent  $i$  adds Gaussian Noise to the prediction (i.e.,  $C$ -dimensional  
136 histogram where each bin is zero except the  $f_i(x_t)$ -th bin is 1). The “pseudo label” is achieved with  
137 the majority vote returned by aggregating the noisy predictions from the local agents.

138 For instance-level DP, the spirit of our method shares with PATE, in the aspect of by adding or  
139 removing one instance, it can *change* at most one agent’s prediction. The same argument also  
140 naturally applies to *adding or removing one agent*. In fact we gain a factor of 2 in the stronger  
141 agent-level DP due to a smaller sensitivity in our approach (see Theorem 4).

142 Another important difference is that in the original PATE, the teacher models are trained on I.I.D data  
143 (random splits of the whole private data), while in our case, the agents are naturally present with  
144 different distributions. We propose to optionally use domain adaptation techniques to mitigate these  
145 differences when training the agents.

### 146 3.2 kNN - DPFL

147 From Definition 2 and 3, preserving agent-level DP is generally more difficult than the instance-  
148 level DP. We find that for *AE-DPFL*, the privacy guarantee for instance-level DP is weaker than its  
149 agent-level DP guarantee (see Theorem 4). To amplify the instance-level DP, we now introduce our  
150 *kNN-DPFL*.

151 As in Algorithm 2, each agent maintains a data-independent feature extractor  $\phi$ , i.e., an ImageNet [Deng et al., 2009]  
152 pre-trained network without the classifier layer. For each unlabeled  
153 query  $x_t$ , agent  $i$  first finds the  $k_i$  nearest neighbors to  $x_t$  from its local data by measuring the  
154 Euclidean distance in the feature space  $\mathcal{R}^{d_\phi}$ . Then,  $f_i(x_t)$  outputs the frequency vector of the votes  
155 from the nearest neighbors, which equals to  $\frac{1}{k}(\sum_{j=1}^k y_j)$ , where  $y_j \in \mathcal{R}^C$  indicates the one-hot  
156 vector of the ground-truth label. Subsequently,  $\tilde{f}_i(x_t)$  from all agents are privately aggregated with  
157 the argmax of the noisy voting scores returned to the server.

158 Our kNN-DPFL differs from Private-kNN in that we apply kNN on each agent’s local data instead of  
159 the entire private dataset. This distinction together with MPC allows us to receive up to  $kN$  neighbors  
160 while bounding the contribution of individual agents by  $k$ . Comparing to *AE-DPFL*, this approach  
161 enjoys a stronger instance-level DP guarantee since the sensitivity from adding or removing one  
162 instance is a factor of  $k/2$  times smaller than that of the agent-level (see the proof in Theorem 4).

163 **How to implement MPC-vote?** Dery et al. [2019] assumes a set of (honest and non-colluding)  
164 external entities, named talliers,  $\mathbb{T} = \{T_1, \dots, T_J\}$ . Then, each agent applies secret sharing for  
165 creating  $J$  shares of the private ballots ( $\tilde{f}_i(x_t)$  in our case), and distributing them among the  $J$  talliers.  
166 After receiving the ballot shares from all agents, the tallier will compute the sum of share vectors  
167 and find the index  $y \in \{1, \dots, C\}$  with the highest scores and send that to the server. We refer the  
168 reader to Protocol 1 in Dery et al. [2019] for a detailed procedure. We highlight that using MPC-vote  
169 (only the top-one index is revealed to the server) instead of MPC-sum results in a stronger differential  
170 privacy guarantee, as discussed in the next section.

### 171 3.3 Privacy Analysis

172 Our privacy analysis is based on Renyi differential privacy (RDP) [Mironov, 2017]. We defer the  
173 background about RDP, its connection to DP and all proofs of our technical results to the appendix  
174 RDP section.

175 **Theorem 4** (Privacy guarantee). *Let AE-DPFL and kNN-DPFL answer  $Q$  queries with noise scale  $\sigma$ .  
176 For agent-level protection, both algorithms guarantee  $(\alpha, \frac{Q\alpha}{2\sigma^2})$ -RDP for all  $\alpha \geq 1$ . For instance-level  
177 protection, AE-DPFL and kNN-DPFL obey  $(\alpha, \frac{Q\alpha}{\sigma^2})$  and  $(\alpha, \frac{Q\alpha}{k\sigma^2})$ -RDP respectively.*

178 **Remark 1.** *Theorem 4 suggests that both algorithms achieve agent-level and instance-level differ-  
179 ential privacy. With the same noise injection to the agent’s output, kNN-DPFL enjoys a stronger  
180 instance-level DP (by a factor of  $k/2$ ) compared to its agent-level guarantee, while AE-DPFL’s  
181 instance-level DP is weaker by a factor of 2. Since AE-DPFL allows an easy-extension with the do-  
182 main adaptation technique, we choose to use AE-DPFL for the agent-level DP and apply kNN-DPFL  
183 for the instance-level DP in the experiments.*

---

**Algorithm 1** *AE-DPFL* with MPC-Vote

---

**input** Noise level  $\sigma$ , unlabeled public data  $\mathcal{D}_G$ , integer  $Q$ .

- 1: Train local model  $f_i$  using  $\mathcal{D}_i$  or using  $(\mathcal{D}_i, \mathcal{D}_G)$  with any domain adaptation techniques.
- 2: **for**  $t = 0, 1, \dots, Q$ , pick  $x_t \in \mathcal{D}_G$  **do**
- 3:   **for** each agent  $i$  in  $1, \dots, N$  (in parallel) **do**
- 4:      $\tilde{f}_i(x_t) = f_i(x_t) + \mathcal{N}(0, \frac{\sigma^2}{N} I_C)$ .
- 5:   **end for**
- 6:    $\tilde{y}_t = \operatorname{argmax}_{y \in \{1, \dots, C\}} [\sum_{i=1}^N \tilde{f}_i(x_t)]_y$  via MPC.
- 7: **end for**

**output** A global model  $\theta$  trained using  $(x_t, \tilde{y}_t)_{t=1}^Q$

---



---

**Algorithm 2** *kNN-DPFL* with MPC-Vote

---

**input** Noise level  $\sigma$ , unlabeled public data  $\mathcal{D}_G$ , integer  $Q$ , feature map  $\phi$ .

- 1: **for**  $t = 0, 1, \dots, Q$ , pick  $x_t \in \mathcal{D}_G$  **do**
- 2:   **for** each agent  $i$  in  $1, \dots, N$  (in parallel) **do**
- 3:     Apply  $\phi$  on  $\mathcal{D}_i$  and  $x_t$
- 4:      $y_1, \dots, y_k \leftarrow$  labels of the  $k$  nearest neighbor.
- 5:      $\tilde{f}_i(x_t) = \frac{1}{k} (\sum_{j=1}^k y_j) + \mathcal{N}(0, \frac{\sigma^2}{N} I_C)$
- 6:   **end for**
- 7:    $\tilde{y}_t = \operatorname{argmax}_{y \in \{1, \dots, C\}} [\sum_{i=1}^N \tilde{f}_i(x_t)]_y$  via MPC.
- 8: **end for**

**output** A global model  $\theta$  trained using  $(x_t, \tilde{y}_t)_{t=1}^Q$

---

184 **Communication Cost:** Finally, we find that our methods are *embarrassingly parallel* as each agent  
 185 work independently without any synchronization. Overall, we reduce the (per-agent) up-stream  
 186 communication cost from  $d \cdot T$  floats (model size times  $T$  rounds) to  $C \cdot Q$ , where  $C$  is number  
 187 of classes and  $Q$  is the number of data points. Moreover, the communication overheads due to  
 188 MPC protocols approach a multiplicative constant over the transmitted data for both MPC-sum and  
 189 MPC-vote ([Bonawitz et al., 2017a, Dery et al., 2019]).

## 190 4 Experimental Results

191 In this section, we apply our *AE-DPFL* for agent-level DP and *kNN-DPFL* for instance-level DP  
 192 based on their distinctive characteristics in privacy guarantee.

### 193 4.1 Agent-level DP Evaluation

194 To investigate various heterogeneous scenarios, we consider: (1) non-I.I.D partition of local data  
 195 (MNIST); (2) data across agents and the server are drawn from different domains (Digit Datasets).

Datasets	# Agents	Methods	Accuracy (%)	$\epsilon$
MNIST (non-I.I.D)	100	FedAvg	$97.8 \pm 0.1$	-
		DP-FedAvg	$84.2 \pm 0.2$	4.3
		<i>AE-DPFL</i> (Ours)	<b><math>86.1 \pm 0.2</math></b>	<b>4.3</b>
SVHN, MNIST	200	FedAvg	$87.6 \pm 0.1$	-
		FedAvg+DA	$86.9 \pm 0.1$	-
		DP-FedAvg	$76.3 \pm 0.3$	3.7
		DP-FedAvg+DA	$71.2 \pm 0.4$	3.6
→ USPS (non-I.I.D)		<i>AE-DPFL</i> (Ours)	$83.8 \pm 0.2$	3.6
		<i>AE-DPFL</i> +DA (Ours)	<b><math>92.5 \pm 0.2</math></b>	<b>2.8</b>

Table 1: **Agent-level DP Evaluation.** We set  $\delta = 10^{-3}$  for all datasets. For MNIST, each local agent is with 6 digits. Different local agents do not share exactly the same 6 digits, which is a non-I.I.D setting. Further, we assign SVHN and MIST for local agents and USPS for the server, which is a typical non-I.I.D with domain shift setting.

196 **MNIST Dataset with Non-I.I.D Partition:** We choose a similar experimental setup as the original  
 197 FedAvg [McMahan et al., 2017] and DP-FedAvg [Geyer et al., 2017] did. We divide the training set  
 198 of the sorted MNIST into 100 agents, such that each agent will have samples from 6 digits only. This  
 199 way, each agent gets 600 data points from 6 classes. We split 30% of the testing set in MNIST as the  
 200 available unlabeled public data and the remaining testing set used for testing.

201 **Digit Datasets Evaluation:** MNIST, SVHN and USPS are put together termed as Digit datasets  
 202 [LeCun et al., 1998, Netzer et al., 2011]. It is a controlled setting to mimic the real situations, where  
 203 distribution of agent-to-server or agent-to-agent can be different. Based on the size of each dataset,

Network	Methods	$A, C, D \rightarrow W$ (Acc. %)	$\epsilon$	$A, C, W \rightarrow D$ (Acc.)	$\epsilon$
AlexNet	FedAvg	$90.5 \pm 0.1$	-	$96.8 \pm 0.1$	-
	DP-FedAvg	$28.1 \pm 0.7$	46.6	$48.2 \pm 0.8$	47.1
	DP-FedSGD	$32.6 \pm 0.9$	4.1	$48.3 \pm 0.9$	4.0
	DP-FedSGD	$75.2 \pm 0.5$	12.4	$83.7 \pm 0.6$	7.9
	<i>kNN-DPFL</i> ( $\sigma = 15$ , Ours)	<b><math>75.4 \pm 0.3</math></b>	<b>3.9</b>	<b><math>84.3 \pm 0.3</math></b>	<b>3.7</b>
ResNet50	FedAvg	$96.5 \pm 0.1$	-	$97.8 \pm 0.1$	-
	DP-FedSGD	$25.8 \pm 0.6$	4.0	$42.7 \pm 0.5$	3.9
	<i>kNN-DPFL</i> ( $\sigma = 25$ , Ours)	<b><math>86.3 \pm 0.4</math></b>	<b>2.8</b>	<b><math>91.9 \pm 0.2</math></b>	<b>2.0</b>

Table 2: **Instance-level DP on Office-Caltech dataset for non-I.I.D setting. Total number of local agents is 3. We set  $\delta = 10^{-4}$ .**

204 we simulate 140 agents using SVHN with 3000 records each and 60 agents using MNIST with 1000  
205 records each. We split 3000 unlabeled records from USPS at server and the rest data is used for  
206 testing.

207 We notice that DP-FedAvg and FedAvg never see the server distribution. To boost those two algo-  
208 rithms, we further apply a standard domain adaptation (DA) technique — adversarial training [Ganin  
209 et al., 2016] on top, denoted as DP-FedAvg+DA and FedAvg+DA, respectively. As a consequence,  
210 their local training involves both local data and unlabeled data from the server. Similarly, we define  
211 *AE-DPFL+DA* as the DA extension of *AE-DPFL*, where each teacher (agent) model is trained with  
212 the same DA technique as that in DP-FedAvg+DA.

213 In Table 1, we observe that when the privacy cost  $\epsilon$  of DP-FedAvg and *AE-DPFL* is close, our method  
214 significantly improves the accuracy from 76.3% to 83.8%. (2) The further improved accuracy 92.5%  
215 of *AE-DPFL+DA* demonstrates that our framework can orthogonally benefit from DA techniques,  
216 where it is highly uncertain yet for the gradient-based methods. (3) Both FedAvg and DP-FedAvg  
217 perform better than their DA variants; therefore we will only use DP-FedAvg in the following  
218 experiments. This result is well expected, as FL with domain adaptation is more closely related to  
219 the multi-source domain adaptation [Peng et al., 2019a]. Combining *FedAvg* with the one-source  
220 DA methods implies averaging different trajectories towards the server’s distribution, which may  
221 not work in practice. Similar learning bound based observation has been investigated in Peng et al.  
222 [2019b] and it remains unclear how to privatize the multi-source domain adaptation approach. On the  
223 other hand, leveraging the majority vote is more stable against the distribution shift. We conjecture  
224 this is because whenever there is a high consensus among the vote counts, the returned label remains  
225 unchanged if the distribution of some agents is slightly perturbed. In contrast, averaging trajectories  
226 in such case may diverge the optimization procedure directly.

## 227 4.2 Instance-level DP Evaluation

228 We investigate the instance-level DP using datasets Office-Caltech10 [Gong et al., 2012] to further  
229 highlight that our method can facilitate to the extreme challenging domain shift scenario, while not  
230 explicitly applying any of the domain adaptation technique. Office-Caltech consists of data from four  
231 domains: Caltech (C), Amazon (A), Webcam(W) and DSLR (D). We iteratively pick one domain  
232 as the server domain each time and the rest ones are for local agents (e.g., in  $A, C, D \rightarrow W$ , W is  
233 treated as the server). For *kNN-DPFL*, we instantiate the public feature extractor using the network  
234 backbone without the classifier layer. The DP-FedSGD method provides the DP baseline where we  
235 use mostly the same parameters as Abadi et al. [2016]. In each experiment, we split 70% data from  
236 the server domain as the public available unlabeled data, which is also the data to be labeled for  
237 *kNN-DPFL*, while the remaining 30% data is used for testing.

238 In Table 2, we observe: (1) DP-FedSGD degrades when the backbone changes from the light  
239 load AlexNet to the heavy load ResNet50, while ours is improved by 10%. It is because larger  
240 model capacity leads to more sensitive response to gradient clipping or noise injection, which has  
241 been surveyed in Abadi et al. [2016]. In contrast, our *kNN-DPFL* avoids the gradient operation  
242 by label aggregation and can still benefit from the larger model capacity. Again, our method  
243 achieves consistently better utility-privacy trade-off as maintaining same privacy cost and can achieve  
244 significantly better utility, or maintaining same utility and can achieve much lower privacy cost.

245 **References**

- 246 Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and  
247 Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*  
248 *Conference on Computer and Communications Security*, pages 308–318, 2016.
- 249 Alekh Agarwal, Martin J Wainwright, Peter L Bartlett, and Pradeep K Ravikumar. Information-  
250 theoretic lower bounds on the oracle complexity of convex optimization. In *Advances in Neural*  
251 *Information Processing Systems*, pages 1–9, 2009.
- 252 Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization: Efficient  
253 algorithms and tight error bounds. In *Proceedings of the 54th Annual IEEE Symposium on*  
254 *Foundations of Computer Science*, pages 464–473, 2014.
- 255 Raef Bassily, Om Thakkar, and Abhradeep Guha Thakurta. Model-agnostic private learning. In  
256 *Advances in Neural Information Processing Systems*, pages 7102–7112, 2018.
- 257 Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan, Sarvar  
258 Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for privacy-  
259 preserving machine learning. In *proceedings of the 2017 ACM SIGSAC Conference on Computer*  
260 *and Communications Security*, pages 1175–1191, 2017a.
- 261 Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan, Sarvar  
262 Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for privacy-  
263 preserving machine learning. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer*  
264 *and Communications Security*, pages 1175–1191, 2017b.
- 265 Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer:  
266 Evaluating and testing unintended memorization in neural networks. In *28th {USENIX} Security*  
267 *Symposium ({USENIX} Security 19)*, pages 267–284, 2019.
- 268 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
269 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
270 pages 248–255. Ieee, 2009.
- 271 Lihi Dery, Tamir Tassa, and Avishay Yanai. Fear not, vote truthfully: Secure multiparty computation  
272 of score based rules. *Expert Systems with Applications*, 168:114434, 2019.
- 273 Irit Dinur and Kobbi Nissim. Revealing information while preserving privacy. In *Proceedings of the*  
274 *twenty-second ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*,  
275 pages 202–210, 2003.
- 276 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in  
277 private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- 278 Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François  
279 Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks.  
280 *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- 281 Robin C Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A client  
282 level perspective. *arXiv preprint arXiv:1712.07557*, 2017.
- 283 Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised  
284 domain adaptation. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages  
285 2066–2073. IEEE, 2012.
- 286 Peter Kairouz, Sewoong Oh, and Pramod Viswanath. The composition theorem for differential  
287 privacy. In *International Conference on Machine Learning (ICML-15)*, 2015.
- 288 Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin  
289 Bhagoji, Keith Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances  
290 and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- 291 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to  
292 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

- 293 Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.  
294 Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.
- 295 Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of  
296 fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019.
- 297 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.  
298 Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelli-*  
299 *gence and Statistics*, pages 1273–1282. PMLR, 2017.
- 300 H. Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private  
301 recurrent language models. In *International Conference on Learning Representations (ICLR-18)*,  
302 2018.
- 303 Ilya Mironov. Rényi differential privacy. In *Computer Security Foundations Symposium (CSF), 2017*  
304 *IEEE 30th*, pages 263–275. IEEE, 2017.
- 305 Payman Mohassel and Yupeng Zhang. Secureml: A system for scalable privacy-preserving machine  
306 learning. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 19–38. IEEE, 2017.
- 307 Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. *arXiv*  
308 *preprint arXiv:1902.00146*, 2019.
- 309 Yurii Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*, volume 87. Springer  
310 Science & Business Media, 2003.
- 311 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading  
312 digits in natural images with unsupervised feature learning. 2011.
- 313 Nicolas Papernot, Martín Abadi, Úlfar Erlingsson, Ian Goodfellow, and Kunal Talwar. Semi-  
314 supervised knowledge transfer for deep learning from private training data. In *International*  
315 *Conference on Learning Representations (ICLR-17)*, 2017.
- 316 Nicolas Papernot, Shuang Song, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and Úlfar Erlings-  
317 son. Scalable private learning with pate. In *International Conference on Learning Representations*  
318 *(ICLR-18)*, 2018.
- 319 Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching  
320 for multi-source domain adaptation. In *Proceedings of the IEEE International Conference on*  
321 *Computer Vision*, pages 1406–1415, 2019a.
- 322 Xingchao Peng, Zijun Huang, Yizhe Zhu, and Kate Saenko. Federated adversarial domain adaptation.  
323 *arXiv preprint arXiv:1911.02054*, 2019b.
- 324 Daniel Peterson, Pallika Kanani, and Virendra J Marathe. Private federated learning with domain  
325 adaptation. *arXiv preprint arXiv:1912.06733*, 2019.
- 326 Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks  
327 against machine learning models. In *IEEE Symposium on Security and Privacy (SP)*, pages 3–18.  
328 IEEE, 2017.
- 329 Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S Talwalkar. Federated multi-task  
330 learning. In *Advances in Neural Information Processing Systems*, pages 4424–4434, 2017.
- 331 Shuang Song, Kamalika Chaudhuri, and Anand D Sarwate. Stochastic gradient descent with differen-  
332 tially private updates. In *Conference on Signal and Information Processing*, 2013.
- 333 Stacey Truex, Nathalie Baracaldo, Ali Anwar, Thomas Steinke, Heiko Ludwig, Rui Zhang, and  
334 Yi Zhou. A hybrid approach to privacy-preserving federated learning. In *AISeC*, 2019.
- 335 Shiqiang Wang, Tiffany Tuor, Theodoros Salonidis, Kin K Leung, Christian Makaya, Ting He, and  
336 Kevin Chan. Adaptive federated learning in resource constrained edge computing systems. *IEEE*  
337 *Journal on Selected Areas in Communications*, 37(6):1205–1221, 2019.

338 Andrew C Yao. Protocols for secure computations. In *23rd annual symposium on foundations of*  
 339 *computer science (sfcs 1982)*, pages 160–164. IEEE, 1982.

340 Yuqing Zhu, Xiang Yu, Manmohan Chandraker, and Yu-Xiang Wang. Private-knn: Practical dif-  
 341 ferential privacy for computer vision. In *The IEEE Conference on Computer Vision and Pattern*  
 342 *Recognition (CVPR)*, June 2020.

## 343 A Challenges in Gradient-based Federated Learning

344 In this section, before introducing our approaches, we motivate them by highlighting the main  
 345 challenges in the conventional DPFL methods in terms of gradient estimation, convergence, and data  
 346 heterogeneity. For other challenges, we refer the readers to a survey [Kairouz et al., 2019].

347 **Challenge 1: Biased Gradient Estimation.** Recent works [Li et al., 2018] have shown that the  
 348 FedAvg may not converge well under data heterogeneity. We provide a simple example to show that  
 349 the clipping step of DP-FedAvg may exacerbate the issue.

350 **Example 5.** Let  $N = 2$ , each agent  $i$ 's local update is  $\Delta_i$  ( $E$  iterations of SGD). We enforce  
 351 clipping of per-agent update  $\Delta_i$  by performing  $\Delta_i / \max(1, \frac{\|\Delta_i\|_2}{S})$ , where  $S$  is the clipping threshold.  
 352 Consider the special case when  $\|\Delta_1\|_2 = S + \alpha$  and  $\|\Delta_2\|_2 \leq S$ . Then the global update will be  
 353  $\frac{1}{2}(\frac{S\Delta_1}{\|\Delta_1\|_2} + \Delta_2)$ , which is biased.

354 Comparing to the FedAvg updates  $\frac{1}{2}(\Delta_1 + \Delta_2)$ , the biased update could be 0 (not moving) or  
 355 pointing towards the opposite direction. Such a simple example can be embedded in more realistic  
 356 problems, causing substantial bias that leads to non-convergence.

357 **Challenge 2: Slow Convergence.** Following works on FL convergence analysis [Li et al., 2019,  
 358 Wang et al., 2019], we derive the convergence analysis on DP-FedAvg and demonstrate that using  
 359 many outer-loop iterations ( $T$ ) could result in similar convergence issue under differential privacy.

360 The appeal of FedAvg is to set  $E$  to be larger so that each agent performs  $E$  iterations to update its  
 361 own parameters before synchronizing the parameters to the global model, hence reducing the number  
 362 of rounds in communication. We show that the effect of increasing  $E$  is essentially increasing the  
 363 learning rate for a large family of optimization problems with piece-wise linear objective functions,  
 364 which does not change the convergence rate. The detailed analysis is in appendix convergence section  
 365 due to space limit. Specifically, it is known that for the family of  $G$ -Lipschitz functions supported  
 366 on a  $B$ -bounded domain, any Krylov-space method <sup>1</sup> has convergence rate that is lower bounded  
 367 by  $\Omega(BG/\sqrt{T})$  [Nesterov, 2003, Section 3.2.1]. This indicates that the variant of FedAvg requires  
 368  $\Omega(1/\alpha^2)$  rounds of outer loop (i.e., communication), in order to converge to an  $\alpha$  stationary point,  
 369 i.e., increasing  $E$  does *not* help, even if no noise is added.

370 It also indicates that DP-FedAvg is essentially the same as *stochastic* sub-gradient method in almost  
 371 all locations of a piece-wise linear objective function with gradient noise being  $\mathcal{N}(0, \sigma^2/N I_d)$ .  
 372 The additional noise in DP-FedAvg imposes more challenges to the convergence. If we plan to

373 run  $T$  rounds and achieve  $(\epsilon, \delta)$ -DP, we need to choose  $\sigma = \frac{\eta EG \sqrt{2T \log(1.25/\delta)}}{N\epsilon}$  [McMahan et al.,  
 374 2018, Theorem 1], which results in a convergence rate upper bound of  $\frac{GB(\sqrt{1 + \frac{2Td \log(1.25/\delta)}{N^2 \epsilon^2}})}{\sqrt{T}} =$   
 375  $O\left(\frac{GB}{\sqrt{T}} + \frac{\sqrt{d \log(1.25/\delta)}}{N\epsilon}\right)$ , for an optimal choice of the learning rate  $E\eta$ .

376 The above bound is tight for stochastic sub-gradient methods, and in fact also information-theoretically  
 377 optimal. The  $GB/\sqrt{T}$  part of upper bound matches the information-theoretical lower bound for all  
 378 methods that have access to  $T$ -calls of stochastic sub-gradient oracle [Agarwal et al., 2009, Theorem  
 379 1]. While the second matches the information-theoretical lower bound for all  $(\epsilon, \delta)$ -differentially  
 380 private methods on the agent level [Bassily et al., 2014, Theorem 5.3]. That is, the first term indicates  
 381 that there must be *many rounds of communications*, while the second term says that the *dependence in*  
 382 *ambient dimension  $d$*  is unavoidable for DP-FedAvg. Clearly, our method also has such dependence  
 383 *in the worst case*. But it is easier for our approach to adapt to the structure that exists in the data  
 384 (i.e., high consensus among voting), as we will illustrate later. In contrast, it has larger impact on

<sup>1</sup>One that outputs a solution in the subspace spanned by a sequence of sub-gradients.

385 DP-FedAvg, since it needs to explicitly add noise with variance  $\Omega(d)$ . Another observation is when  
 386  $N$  is small, no DP method with reasonable  $\epsilon, \delta$  parameters is able to achieve high accuracy for  
 387 agent-level DP. This partially motivates us to consider the other regime that deals with instance-level  
 388 DP.

389 **Challenge 3: Data Heterogeneity.** Federated learning with domain adaptation has been studied  
 390 in Peng et al. [2019b], where they propose a dynamic attention model to adjust the contribution  
 391 from each source (agent) collaboratively. However, most multi-source domain adaptation algorithms,  
 392 including this approach, require sharing local feature vectors to the target domain, which is not  
 393 compatible with the DP setting. Enhancing *DP-FedAvg* with the effective domain adaptation technique  
 394 remains an open problem.

## 395 B Other properties of differential privacy

**Definition 6** (Renyi Differential Privacy [Mironov, 2017]). *We say a randomized algorithm  $\mathcal{M}$  is  $(\alpha, \epsilon(\alpha))$ -RDP with order  $\alpha \geq 1$  if for neighboring datasets  $D, D'$ ,*

$$\mathbb{D}_\alpha(\mathcal{M}(D) || \mathcal{M}(D')) := \frac{1}{\alpha - 1} \log \mathbb{E}_{o \sim \mathcal{M}(D')} \left[ \left( \frac{\Pr[\mathcal{M}(D) = o]}{\Pr[\mathcal{M}(D') = o]} \right)^\alpha \right] \leq \epsilon(\alpha).$$

396 RDP inherits and generalizes the information-theoretical properties of DP.

397 **Lemma 7** (Selected Properties of RDP [Mironov, 2017]). *If  $\mathcal{M}$  obey  $\epsilon_{\mathcal{M}}(\cdot)$ -RDP, then*

1. [Indistinguishability] *For any measurable set  $S \subset \text{Range}(\mathcal{M})$ , and any neighboring  $D, D'$*

$$e^{-\epsilon(\alpha)} \Pr[\mathcal{M}(D') \in S]^{\frac{\alpha}{\alpha-1}} \leq \Pr[\mathcal{M}(D) \in S] \leq e^{\epsilon(\alpha)} \Pr[\mathcal{M}(D') \in S]^{\frac{\alpha}{\alpha-1}}.$$

398 2. [Post-processing] *For all function  $f$ ,  $\epsilon_{f \circ \mathcal{M}}(\cdot) \leq \epsilon_{\mathcal{M}}(\cdot)$ .*

399 3. [Composition]  $\epsilon_{(\mathcal{M}_1, \mathcal{M}_2)}(\cdot) = \epsilon_{\mathcal{M}_1}(\cdot) + \epsilon_{\mathcal{M}_2}(\cdot)$ .

400 This composition rule often allows for tighter calculations of  $(\epsilon, \delta)$ -DP for the composed mechanism  
 401 than the strong composition theorem in [Kairouz et al., 2015]. Moreover, we can convert RDP to  
 402  $(\epsilon, \delta)$ -DP for any  $\delta > 0$  using:

403 **Lemma 8** (From RDP to DP). *If a randomized algorithm  $\mathcal{M}$  satisfies  $(\alpha, \epsilon(\alpha))$ -RDP, then  $\mathcal{M}$  also  
 404 satisfies  $(\epsilon(\alpha) + \frac{\log(1/\delta)}{\alpha-1}, \delta)$ -DP for any  $\delta \in (0, 1)$ .*

405 **Threat models and Multi-Party Computation (MPC)** However, the privacy guarantee of DP-  
 406 FedAvg only applies to the global model and does not apply to the inference made by curious parties  
 407 who can eavesdrop in the network traffics. Cryptographic techniques such as Multi-Party Computation  
 408 (MPC) [Yao, 1982] securely aggregates local updates and ensures privacy against inferences made  
 409 during the communication process. Specifically, if each party adds a small independent noise to  
 410 the part they contribute, MPC ensures that an attacker can only observe the total, even if he taps  
 411 the network messages and hacks into the server. Unfortunately, it is challenging to apply MPC in  
 412 either DP-FedAvg or DP-FedSGD due to high computational overheads. As shown in Bonawitz  
 413 et al. [2017a], the computational cost of security aggregation (used as MPC-Sum in Figure 1) is  
 414  $O(N^2 + dN)$  for users and  $O(dN^2)$  for the server, where  $d$  is the model size and  $N$  is the number of  
 415 agents. In this paper, we consider a new MPC technique due to [Dery et al., 2019] that allows only  
 416 the voted winner to be released while keeping the voting scores completely hidden. This allows us to  
 417 further amplify the DP guarantees. In our experiment, we assume the aggregation is conducted by  
 418 MPC for all privacy-preserving algorithms that we consider (see Figure 1).

## 419 C More Discussions of Challenges for Gradient-Based FL

**Definition 9.** *A function  $\ell$  is Lipschitz continuous with constant  $G > 0$ , if*

$$|\ell(x) - \ell(y)| \leq G \|x - y\|_2$$

420 *for all  $x, y$ .*

421 **Proposition 10.** *Let the objective function of agents  $f_1, \dots, f_N$  obeys that  $f_i$  is piecewise linear*  
422 *(which implies that the global objective  $F = \frac{1}{N} \sum_{i=1}^N f_i$  is piecewise linear) and  $G$ -Lipschitz. Let  $\eta$*   
423 *be the learning rate taken by individual agents. Then the outer loop FedAvg update is equivalent to*  
424  *$\theta^+ = \theta - E\eta g$  for some  $g \in \mathbb{R}^d$ , where (a)  $g = \nabla F(\theta)$  if  $\theta$  is in the  $\nu$  interior of the linear region*  
425 *of  $f_1, \dots, f_N$  and  $E < \nu/(\eta G)$ ; (2)  $g$  is a Clarke-subgradient<sup>2</sup> of  $F$  at  $\theta$ , if  $\theta$  is on the boundary of*  
426 *at least two linear regions and at least  $\nu$  away in Euclidean distance from another boundary and*  
427  *$E < \nu/(\eta G)$ ; (c) otherwise, we have that  $\|g - \nabla F(\theta)\|_2 \leq E\eta G$ . Moreover, statement (c) is true*  
428 *even if we drop the piecewise linear assumption.*

*Proof.* For the Statement (a), observe that for all  $\theta'$  such that  $\|\theta' - \theta\| \leq \nu$  neighborhood, we have that  $\nabla f_i(\theta') = \nabla f_i(\theta)$ . When  $E < \nu/(\eta G)$ , the cumulative gradients of agent  $i$  is equal to  $E\nabla f_i(\theta)$ . For Statement (b), notice that the Clarke subdifferential at  $\theta$  is the convex hull of the one-sided gradient, thus as we move along the negative gradient direction in the inner loop, we enter and remains in the linear region. Thus the update direction is

$$\frac{1}{N} \left( \sum_{i \text{ s.t. } f_i \text{ is differentiable at } \theta} E\eta \nabla f_i(\theta) + \sum_{i \text{ s.t. } f_i \text{ is not differentiable at } \theta} \eta g_i + (E-1)\nabla f_i(\theta - \eta g_i) \right)$$

429 for all  $g_i$  such that it is a Clarke-subgradient of  $f_i$  it can be written as a convex combination. The  
430 proof is complete by observing that the  $1/N \sum_i$  is also a convex combination and by multiplying  
431 and dividing by  $E$ . Statement (c) is a straightforward application of the Lipschitz property which  
432 says that  $E$  steps can at most get you away for  $\eta EG$  and clearly piecewise linear assumption is not  
433 required.  $\square$

434 This proposition says that in almost all  $\theta$ , increasing  $E$  has the effect of increasing the learning  
435 rate of the subgradient “descent” method for piecewise linear objective functions; and increasing  
436 the learning rate of an approximate gradient method in general for Lipschitz objective functions.  
437 It is known that for the family of  $G$ -Lipschitz function supported on a  $B$ -bounded domain, any  
438 Krylov-space method<sup>3</sup> has a rate of convergence that is lower bounded by  $O(BG/\sqrt{T})$  if running for  
439  $T$  iterations. A close inspection of the lower bound construction reveals that the worst-case problem  
440 is  $\min_{\theta \in \mathbb{R}^T} \max_i \theta_i + \|\theta\|^2$ , namely, a regularized piecewise linear function. This is saying that the  
441 variant of FedAvg that aggregates only the loss-function part of the gradient or projects only when  
442 synchronizing essentially requires  $\Omega(1/\alpha^2)$  rounds of outer loop iterations (thus communication) in  
443 order to converge to an  $\alpha$  stationary point, i.e., increasing  $E$  does *not* help, even if no noise is added.

444 **Lemma 11** (Restatement of Lemma ??). *Conditioning on the teachers, for each public data point*  
445  *$x$ , the noise added to each coordinate is drawn from  $\mathcal{N}(0, \sigma^2/N^2)$ , then with probability  $\geq 1 -$*   
446  *$C \exp\{-N^2\gamma(x)^2/8\sigma^2\}$ , the privately released label matches the majority vote without adding*  
447 *noise.*

448 *Proof.* The proof is a straightforward application of Gaussian tail bounds and a union bound over  $C$   
449 coordinates. Specifically,  $\mathbb{P}[Z_{j^*} < -\gamma(x)/2] \leq e^{-\frac{N^2\gamma(x)^2}{8\sigma^2}}$  for the argmax  $j^*$ . For  $j \neq j^*$ ,  $\mathbb{P}[Z_j >$   
450  $\gamma(x)/2] \leq e^{-\frac{N^2\gamma(x)^2}{8\sigma^2}}$ . By a union bound over all coordinates  $C$ , we get that there perturbation from  
451 the boundedness is smaller than  $\gamma(x)/2$ , which implies correct release of the majority votes.  $\square$

452 This lemma implies that for all public data point  $x$  such that  $\gamma(x) \geq \frac{2\sqrt{2\log(C/\delta)}}{N}$ , the output label  
453 matches noiseless majority votes with probability exponentially close to 1.

## 454 D Data-dependent Privacy Analysis

455

<sup>2</sup>Clarke-subgradient is a generalization of the subgradient to non-convex functions. It reduces to the standard (Moreau) subgradient when  $F$  is convex.

<sup>3</sup>One that outputs a solution in the subspace spanned by a sequence of subgradients.

456 **D.1 Privacy Analysis**

457 **Theorem 12** (Restatement of Theorem 4). *Let AE-DPFL and kNN-DPFL answer  $Q$  queries with noise*  
 458 *scale  $\sigma$ . For agent-level protection, both algorithms guarantee  $(\alpha, Q\alpha/(2\sigma^2))$ -RDP for all  $\alpha \geq 1$ .*  
 459 *For instance-level protection, AE-DPFL and kNN-DPFL obey  $(\alpha, Q\alpha/\sigma^2)$  and  $(\alpha, Q\alpha/(k\sigma^2))$ -RDP*  
 460 *respectively.*

461 *Proof.* In AE-DPFL, for query  $x$ , by the independence of the noise added, the noisy sum is identically  
 462 distributed to  $\sum_{i=1}^N f_i(x) + \mathcal{N}(0, \sigma^2)$ . Adding or removing one data instance from will change  
 463  $\sum_{i=1}^N f_i(x)$  by at most  $\sqrt{2}$  in L2. The Gaussian mechanism thus satisfies  $(\alpha, \alpha s^2/2\sigma^2)$ -RDP on the  
 464 instance-level for all  $\alpha \geq 1$  with an L2-sensitivity  $s = \sqrt{2}$ . This is identical to the analysis in the  
 465 original PATE [Papernot et al., 2018].

466 For the agent-level, the L2 and L1 sensitivities are both 1 for adding or removing one agent.

467 In kNN-DPFL, the noisy sum is identically distributed to  $\frac{1}{k} \sum_{i=1}^N \sum_{j=1}^k y_{i,j} + \mathcal{N}(0, \sigma^2)$ . The change  
 468 of adding or removing one agent will change the sum by at most 1, which implies the same L2  
 469 sensitivity and same agent-level protection as AE-DPFL. The L2-sensitivity from adding or removing  
 470 one instance, on the other hand changes the score by at most  $\sqrt{2/k}$  in L2 due to that the instance  
 471 being replaced by another instance, this leads to an improved instance-level DP that reduces  $\epsilon$  by a  
 472 factor of  $\sqrt{\frac{k}{2}}$ .

473 The overall RDP guarantee follows by the composition over  $Q$  queries. The approximate-DP  
 474 guarantee follows from the standard RDP to DP conversion formula  $\epsilon(\alpha) + \frac{\log(1/\delta)}{\alpha-1}$  and optimally  
 475 choosing  $\alpha$ .  $\square$

476 **D.2 Improved accuracy and privacy with large margin**

477 Let  $f_1, \dots, f_N : \mathcal{X} \rightarrow \Delta^{C-1}$  where  $\Delta^{C-1}$  denotes the probability simplex — the soft-label space.  
 478 Note that both algorithms we propose can be viewed as voting of these teachers which outputs a  
 479 probability distribution in  $\Delta^{C-1}$ . First let us define the margin parameter  $\gamma(x)$  which measures the  
 480 difference between the largest and second largest coordinate of  $\frac{1}{N} \sum_{i=1}^N f_i(x)$ .

481 **Lemma 13.** *Conditioning on the teachers, for each public data point  $x$ , the noise added to each*  
 482 *coordinate is drawn from  $\mathcal{N}(0, \sigma^2/N^2)$ , then with probability  $\geq 1 - C \exp\{-N^2\gamma(x)^2/8\sigma^2\}$ , the*  
 483 *privately released label matches the majority vote without adding noise.*

484 *Proof.* The proof is a straightforward application of Gaussian tail bounds and a union bound over  $C$   
 485 coordinates. Specifically,  $\mathbb{P}[Z_{j^*} < -\gamma(x)/2] \leq e^{-\frac{N^2\gamma(x)^2}{8\sigma^2}}$  for the  $\text{argmax } j^*$ . For  $j \neq j^*$ ,  $\mathbb{P}[Z_j >$   
 486  $\gamma(x)/2] \leq e^{-\frac{N^2\gamma(x)^2}{8\sigma^2}}$ . By a union bound over all coordinates  $C$ , we get that there perturbation from  
 487 the boundedness is smaller than  $\gamma(x)/2$ , which implies correct release of the majority votes.  $\square$

488 This lemma implies that for all public data point  $x$  such that  $\gamma(x) \geq \frac{2\sqrt{2\log(C/\delta)}}{N}$ , the output label  
 489 matches noiseless majority votes with probability exponentially close to 1.

490 Next we show that for those data point  $x$  such that  $\gamma(x)$  is large, the privacy loss for releasing  
 491  $\text{argmax}_j [\frac{1}{N} \sum_{i=1}^N f_i(x)]_j$  is exponentially smaller. The result is based on the following privacy  
 492 amplification lemma that is a simplification of Theorem 6 in the appendix of [Papernot et al., 2018].

**Lemma 14.** *Let  $\mathcal{M}$  satisfy  $(2\alpha, \epsilon)$ -RDP, and there is a singleton output that happens with probability  
 $1 - q$  when  $\mathcal{M}$  is applied to  $D$ . Then for any  $D'$  that is adjacent to  $D$ , Renyi-divergence*

$$D_\alpha(\mathcal{M}(D) \parallel \mathcal{M}(D')) \leq -\log(1 - q) + \frac{1}{\alpha - 1} \log(1 + q^{1/2}(1 - q)^{\alpha-1} e^{(\alpha-1)\epsilon}).$$

493 *Proof.* Let  $P, Q$  be the distribution of  $\mathcal{M}(D)$  and  $\mathcal{M}(D')$  respectively and  $E$  be the event that the  
 494 singleton output is selected.

$$\begin{aligned} \mathbb{E}_Q[(dP/dQ)^\alpha] &= \mathbb{E}_Q[(dP/dQ)^\alpha | E] \mathbb{P}_Q[E] + \mathbb{E}_Q[(dP/dQ)^\alpha \mathbf{1}(E^c)] \\ &\leq (1-q) \left(\frac{1}{1-q}\right)^\alpha + \sqrt{\mathbb{E}_Q[(dP/dQ)^{2\alpha}]} \sqrt{\mathbb{E}_Q[\mathbf{1}(E^c)^2]} \\ &\leq (1-q)^{-(\alpha-1)} + q^{1/2} e^{(2\alpha-1)\epsilon/2} = (1-q)^{-(\alpha-1)} \left(1 + (1-q)^{\alpha-1} q^{1/2} e^{\frac{2\alpha-1}{2}\epsilon}\right) \end{aligned}$$

495 The first part of the second line uses the fact that event  $E$  is a singleton with probability larger than  
 496  $1-q$  under  $Q$  and the probability is always smaller than 1 under  $P$ . The second part of the second  
 497 line follows from Cauchy-Schwartz inequality. The third line substitute the definition of  $(2\alpha, \epsilon)$ -RDP.  
 498 Finally, the stated result follows by the definition of the Renyi divergence.  $\square$

**Theorem 15** (Restatement of Theorem ??). *The mechanism that releases  $\operatorname{argmax}_j [\frac{1}{N} \sum_{i=1}^N f_i(x) + \mathcal{N}(0, (\sigma^2/N^2)I_C)]_j$  obeys  $(\alpha, \epsilon)$ -data-dependent-RDP, where*

$$\epsilon \leq 2C e^{-\frac{N^2 \gamma(x)^2}{8\sigma^2}} + \frac{1}{\alpha-1} \log \left( 1 + e^{\frac{(2\alpha-1)\alpha s}{2\sigma^2} - \frac{N^2 \gamma(x)^2}{8\sigma^2} + \log C/2} \right),$$

499 where  $s = 1$  for AE-DPFL with the agent-level DP, and  $s = 2/k$  for KNN-DPFL with the instance-  
 500 level DP.

501 *Proof.* The proof involves substituting  $q = C e^{-\frac{N^2 \gamma(x)^2}{8\sigma^2}}$  from Lemma ?? into Lemma 14 and use the  
 502 fact that  $\mathcal{M}$  satisfies the RDP of a Gaussian mechanism from the RDP's post-processing lemma. The  
 503 expression bound is simplified for readability using  $-\log(1-x) < 2x$  for all  $x > -0.5$  and that  
 504  $(1-q)^{\alpha-1} \leq 1$ .  $\square$

505 As we can see, when given teachers that are largely in consensus, the (data-dependent) privacy loss  
 506 exponentially smaller.