
HYDRA-FL: Hybrid Knowledge Distillation for Robust and Accurate Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Data heterogeneity among Federated Learning (FL) users poses a significant chal-
2 lenge, resulting in reduced global model performance. The community has de-
3 signed various techniques to tackle this issue, among which Knowledge Distillation
4 (KD)-based techniques are common. While these techniques effectively improve
5 performance under high heterogeneity, they inadvertently cause higher accuracy
6 degradation under model poisoning attacks (known as *attack amplification*). This
7 paper presents a case study to reveal this critical vulnerability in KD-based FL
8 systems. We show why KD causes this issue through empirical evidence and use it
9 as motivation to design a hybrid distillation technique. We introduce a novel algo-
10 rithm, *Hybrid Knowledge Distillation for Robust and Accurate FL (HYDRA-FL)*,
11 ¹, which reduces the impact of attacks in attack scenarios by offloading some of
12 the KD loss to a shallow layer via an auxiliary classifier. We model HYDRA-FL
13 as a generic framework and adapt it to two KD-based FL algorithms, FedNTD
14 and MOON. Using these two as case studies, we demonstrate that our technique
15 outperforms baselines in attack settings while maintaining comparable performance
16 in benign settings.

17 1 Introduction

18 Federated Learning (FL) [32] is an emerging machine learning paradigm enabling multiple users’
19 collaborative model training without data sharing. Each user, termed a *client*, only shares their local
20 model with a *server*, which aggregates all local models into a single global model and redistributes
21 it to the clients. Due to its decentralized, privacy-preserving, and highly-scalable nature, FL has
22 been adopted by Google’s Gboard [2] for next-word prediction, Apple’s Siri [1] for automatic speech
23 recognition, and WeBank [43] for credit risk prediction.

24 Despite its benefits, FL faces challenges with data heterogeneity [28, 51, 13, 24]. FL performs well
25 when client data is independent and identically distributed (IID) and achieves similar convergence
26 as a single model trained on all the clients’ data but struggles when clients have diverse data (non-
27 IID). In this case, the client’s local data is not a good representation of the overall data distribution
28 (unlike an ideal IID case), causing local models to *drift* away from each other. This drift results
29 in a global model with significant accuracy degradation compared to the IID scenario. Numerous
30 solutions [25, 20, 22, 53, 23, 46, 15, 27] address data heterogeneity, including Knowledge Distillation
31 (KD) [12] to reduce the drift between local models.

32 Besides data heterogeneity, FL also faces the issue of Byzantine robustness [14], where *untrusted*
33 *clients* can inject *poisoned* models into the aggregator by altering client data (data poisoning [35])
34 or client models (model poisoning [11, 4, 33, 45, 5, 3, 41]). Research by [40] shows that model

¹We will release the open source code with the final version of this paper.

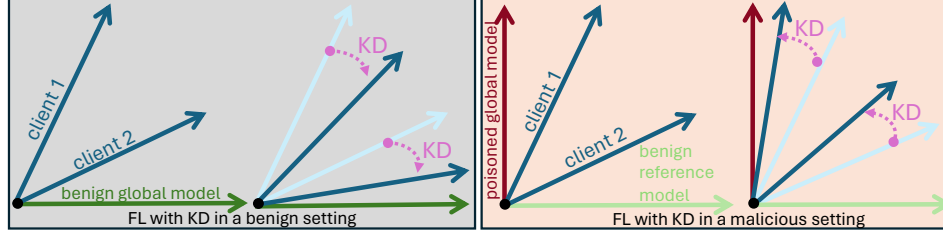


Figure 1: Overview of attack amplification through knowledge distillation. **a)** In the benign setting, KD reduces drift and brings benign local models closer to the benign global model. **b)** In the malicious setting, KD *unknowingly* reduces drift between benign local models and the poisoned global model.

poisoning attacks are more potent as they directly manipulate local models. To counteract poisoning in FL, various defenses have been developed [6, 47, 50, 7, 26, 9, 8].

In this work, we identify a critical vulnerability in KD-based FL techniques under model poisoning attacks. These techniques unknowingly align benign client models with a poisoned server model (Figure 1). We study two such classes of KD-based solutions: FedNTD [20], which reduces the loss between not-true logits of the server and client models, and MOON [25], which reduces the contrastive loss between the representation vector of the server and client models. While these techniques improve global model accuracy in benign settings compared to FedAvg [32] (standard FL aggregator), they *reduce performance below FedAvg under attack*, a phenomenon we term *attack amplification*, especially noticeable at higher heterogeneity levels.

Motivated by our findings, we propose a Hybrid Knowledge Distillation for Robust and Accurate FL (HYDRA-FL) framework for KD-based techniques that restricts attack amplification under poisoning attacks while retaining performance in the benign setting. Unlike traditional KD methods that apply KD-loss only at the final layer, HYDRA-FL introduces KD-loss at a *shallow layer* via an auxiliary classifier and reduces the KD-loss impact at the final layer. This approach draws inspiration from Self-Distillation (SD) [49] and Skeptical Students (SS) [18], but with a distinct focus on enhancing robustness against heterogeneity and model poisoning attacks in FL. SD improves model accuracy by self-distillation, while SS distills from "nasty teachers" [30] to shallow layers. In contrast, our approach uses auxiliary classifiers to enhance FL client robustness against heterogeneity and model poisoning attacks. We design a generic loss function adaptable to specific KD-based algorithms. Extensive experiments show that HYDRA-FL significantly boosts accuracy over FedNTD and MOON in attack settings while maintaining performance in benign settings.

Contributions. This work addresses the critical issue of attack amplification in KD-based FL techniques to counter data heterogeneity. In doing so we make the following contributions:

- **Proving KD amplifies model poisoning:** our motivational case study (§3) on two KD-based techniques, FedNTD and MOON, shows that KD improves accuracy in benign settings but helps the malicious clients propagate poisoning through the KD-loss in adversarial settings. We empirically and theoretically show that this attack amplification issue is inherent to any technique aligning client outputs/representations with the server.
- **Designing HYDRA-FL:** Using our observations as a guideline, we design HYDRA-FL (§4) to prevent attack amplification while retaining performance in the benign setting. HYDRA-FL is formulated as a general loss function adaptable to any FL algorithm to use as its local model training objective.
- **Implementation and Evaluation:** we adapt HYDRA-FL to FedNTD and MOON and modify their local training objectives (§5). Our qualitative and quantitative analysis (§6) shows HYDRA-FL achieves higher accuracy in attack settings and maintains accuracy in benign settings.

72 2 Background and Related Work

73 2.1 Federated Learning (FL)

74 In FL [14, 32], a service provider, called *server*, trains a *global model*, θ^g , on the private data from
 75 multiple collaborating clients, all without directly collecting their data. The server selects n out of
 76 total N clients in every FL round and shares the most recent global model (θ_g^t) with them, where t is
 77 the round number. Then, a client k uses their local data D_k to compute an update ∇_k^t and shares it
 78 with the server. The server aggregates these updates using some *aggregation rule*, like FedAvg [32]
 79 algorithm. In *FedAvg*, a client k *fine-tunes* θ_g^t on their local data using stochastic gradient descent
 80 (SGD) for a fixed number of local epochs E , resulting in an updated local model θ_k^t . The client then
 81 computes their update as the difference $\nabla_k^t = \theta_k^t - \theta_g^t$ and shares ∇_k^t with the server. Next, the server
 82 computes an aggregate of client updates, f_{agg} using mean, i.e.,

$$\nabla_{\text{agg}}^t = f_{\text{mean}}(\nabla_{\{k \in [n]\}}^t). \quad (1)$$

83 The server then updates the global model of the $(t + 1)^{th}$ round using SGD and server learning η as:

$$\theta_g^{t+1} \leftarrow \theta_g^t + \eta \nabla_{\text{agg}}^t \quad (2)$$

84 2.1.1 Data Heterogeneity in FL

85 Data heterogeneity is a well-explored problem [28, 51, 13, 24] in FL. Each client in FL generates
 86 its data, leading to local data distributions that vary across clients and do not accurately represent
 87 the global data distribution. By extension, a global model learned by aggregating local models
 88 using FedAvg may not be the best representation of all the client’s local data. Studies have shown
 89 that this data heterogeneity degrades performance and have proposed various methods to address
 90 this issue [25, 20, 22, 53, 23, 46, 15, 27]. This degradation is more prominent in the presence of
 91 poisoning attacks. Research on poisoning attacks in FL has demonstrated that such attacks become
 92 more successful under high heterogeneity [11, 41]. This increased risk is because the malicious
 93 clients can more easily hide between drifted benign client models, making it difficult for the server
 94 to differentiate between heterogeneous benign clients and malicious ones. [16] highlights that
 95 overlooking this heterogeneity is a critical oversight in FL defense evaluations.

96 2.1.2 Poisoning in FL

97 FL is vulnerable to poisoning attacks [6, 4, 5, 3, 33, 11, 31, 45, 35, 41], where malicious clients aim
 98 to compromise the training process by degrading the global model’s performance. These attacks
 99 come in various forms: In *data poisoning* [3], malicious clients poison their local data to introduce a
 100 backdoor in the local model. This backdoor then propagates to the global model upon aggregation.
 101 In *model poisoning* [11, 4, 33, 45, 5, 3, 41], malicious clients perturb their local models so that,
 102 when aggregated, the global model is poisoned. Poisoning attacks can be further classified based on
 103 their targets: If the performance degradation is on specific inputs, the attack is termed as *targeted*
 104 *poisoning* [5, 3], and if it is on all inputs, then it is termed as *untargeted poisoning* [11, 4, 33, 31, 45].
 105 We explain the attacks used in this paper in §C.2.

106 2.2 Knowledge Distillation (KD)

107 Knowledge Distillation (KD) [12] transfers knowledge from a large, complex model (*teacher*) to
 108 a smaller, more computationally efficient model (*student*). This process involves distilling the
 109 teacher’s rich and intricate information into the student by aligning their predictions. Formally,
 110 if the teacher and student models produce the output probabilities y_t^i and y_s^i respectively for the
 111 i^{th} input (x^i, y^i) , KD aims to match these probabilities by applying the Kullback-Leibler (KL)
 112 divergence between them. The KL-divergence between their softened probabilities is given by:
 113 $KL(\text{softmax}(y_t^i/\tau) || \text{softmax}(y_s^i/\tau))$, where τ is the temperature parameter that softens the proba-
 114 bilities. The overall KD loss function combines this KL-divergence with the usual loss function such
 115 as cross-entropy (CE) loss with β (balances the importance of the KL-divergence and CE loss) as:

$$\mathcal{L} = (1 - \beta) \cdot \mathcal{L}_{CE}(y_s^i, y^i) + \beta \cdot \mathcal{L}_{KL}(\text{softmax}(y_s^i/\tau) || \text{softmax}(y_t^i/\tau)) \quad (3)$$

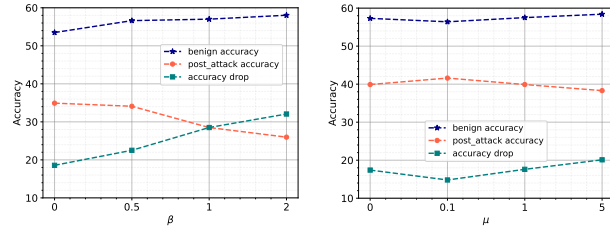
116 **KD in FL** is becoming essential as it addresses critical challenges such as non-IID data distributions,
 117 enhances model performance, accelerates convergence, reduces communication overhead, and im-
 118 proves robustness by making the global model learn from an ensemble of local models [10, 22, 29, 52].
 119 In FL, data is often non-IID across clients, leading to significant discrepancies in local models. KD
 120 mitigates these discrepancies by aligning the local models with the global model, ensuring that
 121 the global model captures a more generalized representation of the data. The general approach is
 122 to reduce the local model drift by improving the aggregation through distillation using unlabeled
 123 auxiliary data. However, the auxiliary data may not always be available, and methods have also been
 124 developed to enable KD without such data [48, 53].

125 3 Attack Amplification through Knowledge Distillation

126 **Hypothesis.** KD-based techniques in FL improve accuracy in non-adversarial settings but result
 127 in more significant accuracy degradation under model poisoning attacks compared to the baseline
 128 techniques such as FedAvg.

129 **Motivational case study.** In this case study, we compare FedAvg against two distinct KD-based
 130 solutions addressing the local model drift from non-IID. MOON [25] uses model-contrastive learning
 131 to align local and global model *representations*, while FedNTD [20] uses KL-divergence to align
 132 *not-true logits* of client models with those of the server. FedNTD penalizes prediction divergence
 133 measured through distillation loss, improving knowledge transfer and stability, while MOON pe-
 134 nalizes *representation divergence* measured through contrastive loss, enhancing robustness and
 135 generalization. This comparison will help us understand the trade-offs of using KD in FL, especially
 136 under adversarial conditions. Throughout this paper, benign conditions mean that no attacks are
 137 present, while adversarial conditions mean that model poisoning attacks are present. We implement
 138 the same settings and hyperparameters for FedAvg as for MOON and FedNTD to ensure a fair
 139 comparison, so FedAvg results may vary between these techniques. This is not an inconsistency. *We*
 140 *do not directly compare FedNTD to MOON unless stated otherwise*, as the original FedNTD work
 141 already did so. Our goal is to test how adversarial settings affect these two fundamentally different
 142 techniques similarly, demonstrating that *attack amplification is inherent to KD and not specific to a*
 143 *particular technique*.

144 **Adversarial conditions.** We simulate
 145 untargeted model poisoning attacks
 146 using techniques from [41, 11]. To
 147 observe their effects on accuracy in
 148 both benign and adversarial settings,
 149 we vary key hyperparameters — LD-
 150 divergence loss coefficient β for Fed-
 151 NTD and contrastive loss coefficient
 152 μ for MOON. The baseline for com-
 153 parison is FedAvg with $\beta = 0$ and
 154 $\mu = 0$. To ensure high heterogeneity
 155 in both settings, the Dirichlet distribu-
 156 tion [34] parameter α is fixed at 0.1.



(a) FedNTD, $\beta = 0$ is FedAvg (b) MOON, $\mu = 0$ is FedAvg

Figure 2: Impact of increasing KL-divergence loss for FedNTD and contrastive loss for MOON on accuracy.

157 **Findings.** In Figures 2(a) and 2(b), we present three key results: benign accuracy (blue), post-
 158 attack accuracy (orange), and the accuracy drop (green). We make the following observations from
 159 increasing β and μ are as follows: (1) the global model accuracy improves in benign settings; (2)
 160 post-attack accuracy decreases; and (3) accuracy drop increases. Our analysis shows a significant
 161 trade-off: *the very mechanisms that improve performance in benign conditions (increasing β and μ)*
 162 *also make the models more vulnerable to adversarial attacks*.

163 **What causes attack amplification?** The fundamental nature of KD-based FL methods aims to align
 164 local models with the global model. In benign scenarios, these methods significantly outperform
 165 FedAvg [25, 20]. However, in the presence of model poisoning attacks, *this model alignment process*
 166 *inadvertently forces local models to align its representation/predictions to the poisoned global model,*
 167 *amplifying the attack's impact*. This is illustrated in Figure 1, where clients *unknowingly distill*
 168 *knowledge* from a poisoned server model.

169 **Formally:** Consider a set of n clients c_1, c_2, \dots, c_n with m being malicious. Using an aggregation
 170 rule such as FedAvg, the server aggregates updates from both benign ($\nabla_{i \in [m+1, n]}$) and malicious
 171 ($\nabla_{i \in [m]}^m$) clients:

$$\nabla_g = f_{\text{agr}}(\nabla_{i \in [m]}^m \cup \nabla_{i \in [m+1, n]}) \quad (4)$$

172 When $m = 0$, the server model ∇_g^b is benign. For $m \neq 0$, the server model ∇_g' is poisoned, deviating
 173 from the ideal unpoisoned global model due to the nature of these attacks [41, 11, 40]. Aligning local
 174 models with a poisoned global model reduces gradient diversity, making local models more similar
 175 to the poisoned global model [20] through KL-divergence or contrastive loss. We rewrite Equation 3
 176 to formalize the loss function for an FL client, using KD, where the client is the student with output
 177 \hat{y}_c , and the server is the teacher with output y_s :

$$\mathcal{L} = \mathcal{L}_{CE}(\hat{y}_c, y) + \beta \mathcal{L}_{KL}(\hat{y}_c, y_s) \quad (5)$$

178 Note that for the sake of derivation here, we are using \hat{y}_c , which represents the generic client model
 179 output. In the case of FedNTD, it can be replaced by \hat{y}_c that represents the not-true logits of the client
 180 model, and in the case of MOON, it can be replaced by z_c that represents the client model's high
 181 dimensional representation.

182 In benign scenarios, this loss function ($\mathcal{L} = \mathcal{F}(\beta)$) decreases monotonically with β because KD
 183 brings local models closer to an unpoisoned global model. Conversely, in adversarial scenarios, it
 184 increases with β because KD brings local models to the poisoned global model. We can write the
 185 relation of this loss function with β as:

$$\mathcal{L}(\beta) \text{ is } \begin{cases} \text{monotonically decreasing,} & m = 0 \\ \text{monotonically increasing,} & m \neq 0 \end{cases} \quad (6)$$

186 Then, the derivative of the loss function is:

$$\frac{d\mathcal{L}}{d\beta} = \begin{cases} < 0, & m = 0 \\ > 0, & m \neq 0 \end{cases} \quad (7)$$

187 Our derivation shows that while the distillation process decreases loss in the absence of malicious
 188 clients, it increases loss in their presence, thereby leading to reduced global model accuracy. This
 189 formal analysis highlights the need for a solution that mitigates the accuracy degradation under
 190 adversarial conditions while retaining the benefits of KD under benign conditions.

191 Impact of Heterogeneity.

192 Now, we explore the effect of het-
 193 erogeneity on the performance
 194 of FedNTD, MOON, and Fe-
 195 dAvg in both benign and adver-
 196 sarial conditions to gain deeper
 197 insights into the role of hetero-
 198 geneity in the KD performance
 199 gain vs. vulnerability tradeoff.
 200 As shown in Figure 3(a), several
 201 interesting observations emerge.
 202 First, both FedNTD and FedAvg
 203 achieve higher accuracy at lower
 204 heterogeneity levels (indicated
 205 by higher α). In benign settings,
 206 FedNTD consistently outperforms FedAvg. However, *the trend reverses in adversarial settings:*
 207 FedAvg achieves higher accuracy than FedNTD, except at $\alpha = 0.5$. A similar pattern is observed
 208 with MOON in Figure 3(b), where FedAvg outperforms MOON across all heterogeneity levels in
 209 adversarial settings. In the benign setting, as expected, MOON slightly outperforms FedAvg at high
 210 heterogeneity. This comparison highlights again how the alignment mechanisms in FedNTD and
 211 MOON with higher heterogeneity exacerbate the vulnerability of KD methods to attacks.

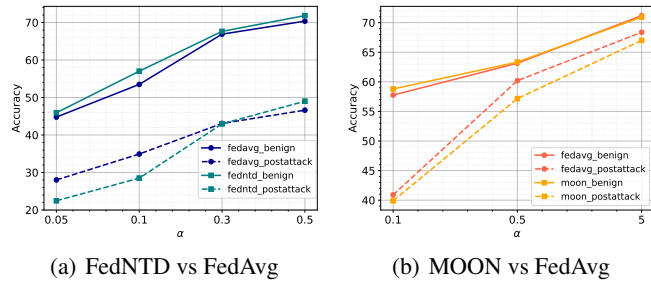


Figure 3: Impact of the heterogeneity parameter, α in benign and adversarial settings. We use the Dirichlet distribution where a higher α means lower heterogeneity.

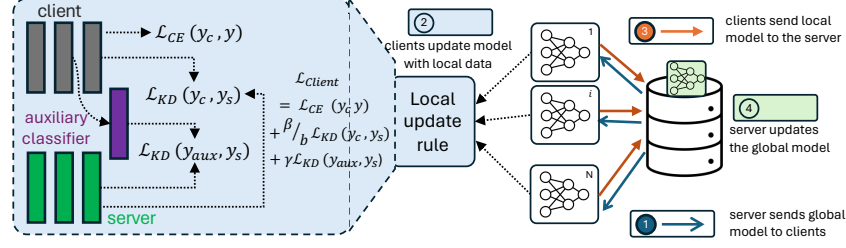


Figure 4: HYDRA-FL framework: we refine client model training by reducing the final layer’s KD-loss and incorporating shallow KD-loss at an earlier shallow layer via an auxiliary classifier.

4 HYDRA-FL: Hybrid Knowledge Distillation for Robust and Accurate FL

4.1 Generic Formulation

In this section, we propose Hybrid Knowledge Distillation for Robust and Accurate FL (HYDRA-FL), a technique to mitigate the *attack amplification* caused by KD in FL. We take a hybrid distillation approach, applying KD-loss at both the final and a shallow layer of the client model (Figure 4). This method incorporates shallow distillation, which applies KD-loss at an intermediate layer and helps reduce the impact of poisoning by preventing over-reliance on final layer alignment. Shallow distillation previously used to handle *nasty teachers* trained adversarially [18], to reduce the impact of poisoning. In summary, shallow layers capture basic features, and shallow distillation ensures these features are robustly learned, protecting the model from adversarial influences that could corrupt deeper layers and final outputs. We first formulate the generic loss function of an FL client using KD in HYDRA-FL as:

$$\mathcal{L} = \mathcal{L}_{CE}(y_c, y) + \frac{\beta}{b} \mathcal{L}_{KD}(y_c, y_s) + \gamma \mathcal{L}_{KD}(y_{aux}, y_s) \quad (8)$$

This loss function has three key components:

1. **Cross-entropy loss** ($\mathcal{L}_{CE}(y_c, y)$) is the loss between the client’s prediction y_c and the target y , drives the client model to learn from its own data, ensuring it captures *in-distribution knowledge* such as features and patterns specific to its data.
2. **Diminished KD loss** ($\frac{\beta}{b} \mathcal{L}_{KD}(y_c, y_s)$) is the loss between the client’s output/representation y_c and the server’s output/representation y_s ². It is a strategic reduction of the KD loss to ensure that the local model benefits from the global model’s knowledge while remaining robust against adversarial attacks. This approach helps balance the trade-offs between learning efficiency and model integrity. In practice, this is achieved by introducing a *diminishing factor* b to the KD loss at the client model’s output layer to diminish the poisoning effect. The KD loss coefficient β is divided by b , effectively reducing its weight in the total loss calculation, thus reducing its influence on the local model’s training. This diminishing factor is essential, as shown later in §6.2 and Figure 7, where reducing the β yields better results.
3. **Shallow distillation loss** ($\gamma \mathcal{L}_{KD}(y_{aux}, y_s)$) is applied at a shallow layer of the local model, enhancing robustness without heavily relying on the final layer alignment. This loss, between the auxiliary classifier’s output/representation y_{aux} at the client model’s shallow layer and the server’s output/representation y_s , is scaled by γ to control the amount of distillation. This approach reduces the impact of poisoning on the client model. Simply reducing the KD-loss in FedNTD or MOON improves post-attack accuracy but reduces benign setting accuracy, as shown in Figure 2. Our shallow distillation loss helps maintain the balance between accuracy in benign settings and lowering the impact of poisoning on the client model in adversarial settings.

Differences with previous works. The key difference between our work and [18] lies in our approach to shallow distillation. [18] aims to distill from models that are designed to be undistillable, a.k.a

² y_c and y_s in generic \mathcal{L}_{KD} loss can be either outputs or representations, because the method can involve either type of comparison (e.g., MOON uses representation-based loss while FedNTD has output-based loss.)

249 *nasty teachers* [30]. While both use hybrid shallow distillation, [18] completely removes the KD-loss
 250 from the model’s output layer and uses self-distillation to compensate for performance loss due to
 251 shallow distillation. In contrast, we retain a scaled-down KD-loss at the output layer. We found that
 252 completely removing the KD-loss at the output layer may cause a more negative impact than keeping
 253 it in a reduced form. Additionally, the untargeted poisoning is different from the poisoning in the
 254 "nasty teacher" paper [30]. The "nasty teacher" performs near-perfect under normal conditions unless
 255 a malicious model distills from it. In untargeted FL poisoning, the global model is poisoned and
 256 performs poorly regardless of its use for distillation.

257 In HYDRA-FL, we use both final layer and shallow layer distillation to enhance robustness. *Final*
 258 *layer distillation* aligns client outputs with server outputs for consistent predictions, whereas the
 259 *shallow layer distillation* aligns intermediate representations to improve robustness against attacks.
 260 This dual approach reduces vulnerability to poisoning attacks, enhances learning by leveraging
 261 knowledge transfer from multiple layers, and maintains high accuracy in benign settings while being
 262 resilient under attack conditions.

263 4.2 Adapting HYDRA-FL to State-of-Art Techniques

264 In this section, we will adapt our generic HYDRA-FL to two state-of-the-art KD techniques for FL.

265 **FedNTD with shallow distillation and auxiliary classifiers.** We modify the FedNTD base model
 266 by introducing auxiliary classifiers. The base model includes two convolutional layers, a linear layer,
 267 and a classification layer. Auxiliary classifiers, each consisting of a linear layer (hidden dimension
 268 512) followed by a classification layer, are added after each convolutional layer. We update the loss
 269 function to include a shallow-distillation term, representing the KL-divergence loss between the
 270 not-true logits of an auxiliary classifier and the global model. The final loss function is a weighted
 271 sum of the standard cross-entropy loss, KL-divergence loss between the not-true logits of the global
 272 model and the client model, and the KL-divergence loss between the not-true logits of the global
 273 model and the auxiliary classifier. The revised loss function in Equation 8 for FedNTD is:

$$\mathcal{L} = \mathcal{L}_{CE}(y_c, y) + \frac{\beta}{b} \mathcal{L}_{KL}(\tilde{y}_c, \tilde{y}_s) + \gamma \mathcal{L}_{KL}(\tilde{y}_{aux}, \tilde{y}_s) \quad (9)$$

274 Here y is the target label, y_c is the client model’s output, \tilde{y}_s , \tilde{y}_c , and \tilde{y}_{aux} are the client model’s,
 275 server model’s, and auxiliary classifier’s not-true logits respectively.

276 **MOON with shallow distillation and auxiliary classifiers.** MOON base model has two convolution
 277 layers, two linear layers, and an output classification layer. We insert auxiliary classifiers after each
 278 convolution layer. Each auxiliary classifier has two linear layers, with a hidden dimension of 256 and
 279 an output dimension of 10. We adapt Equation 8 to MOON to compute the contrastive loss at the
 280 hidden representation layer of the auxiliary classifier as:

$$\mathcal{L} = \mathcal{L}_{CE}(y_c, y) + \frac{\mu}{b} \mathcal{L}_{con}(z_c, z_s) + \gamma \mathcal{L}_{con}(z_{aux}, z_s) \quad (10)$$

281 Here y is the target label, y_c is the client’s output, z_c is the representation from the client’s final
 282 layer, z_s is the representation from the server’s final layer, z_{aux} is the representation from the client’s
 283 auxiliary classifier, and y_s is the server model’s output. For simplicity, we do not write the previous
 284 round’s representation in the loss function here.

285 5 Experimental Results

286 5.1 Experimental Settings

287 **Datasets and Models:** We conduct our experiments over three popular datasets: MNIST, CIFAR10,
 288 and CIFAR100. To ensure a fair comparison with previous works, MOON and FedNTD, we utilized
 289 the same models and hyperparameters they used. Specifically, we incorporated our algorithm as a
 290 simple modification into their publicly available codes [21, 37] (more details in Appendix D).

291 5.2 Shallow Not-True Distillation

292 Our hybrid shallow not-true distillation technique significantly improves post-attack accuracy over
 293 the baseline FedNTD. As shown in Table 1, we achieve higher post-attack accuracy across all

Table 1: Test accuracy for three techniques on three datasets. In the no-attack setting, ($\uparrow\downarrow$) shows comparison to FedAvg. In the attack setting, we use bold if our technique outperforms FedNTD.

Dataset	MNIST		CIFAR10						CIFAR100	
			$\alpha = 0.05$		$\alpha = 0.1$		$\alpha = 0.5$			
Techniques	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>
Fedavg	92.12	74.48	44.69	31.27	54.67	35.67	70.57	48.27	26.17	12.92
FedNTD	93.03 \uparrow	58.09	46.94 \uparrow	21.72	56.95 \uparrow	32.61	71.79 \uparrow	52.51	29.1 \uparrow	13.92
HYDRA-FL(Ours)	92.69 \uparrow	76.67	46.92 \uparrow	25.15	57.12 \uparrow	34.25	71.22 \uparrow	52.57	28.9 \uparrow	14.33

Table 2: Test accuracy for three techniques on three datasets. In the no-attack setting, ($\uparrow\downarrow$) shows comparison to FedAvg. In the attack setting, we use bold if our technique outperforms MOON.

Dataset	MNIST		CIFAR10						CIFAR100	
			$\alpha = 0.1$		$\alpha = 0.5$		$\alpha = 5$			
Methods	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>
Fedavg	88.02	77.55	57.76	40.9	63.14	60.2	71.19	68.38	28.36	24.21
MOON	91.13 \uparrow	72.32	58.8 \uparrow	39.9	63.34 \uparrow	57.17	70.95 \downarrow	67	29.34 \uparrow	23.81
HYDRA-FL(Ours)	92.04 \uparrow	76.65	60.1 \uparrow	43.6	63.32 \uparrow	59.93	70.55 \downarrow	68.4	29.48 \uparrow	25.18

heterogeneity levels. By retaining a diminished NTD loss at the output layer, we maintain similar accuracy to FedNTD in no-attack scenarios and, in some cases, even achieve slightly higher accuracy. We also compare no-attack and post-attack accuracies for FedAvg, the foundational algorithm for many FL aggregation methods.

5.3 Shallow MOON

Our shallow-distillation design effectively prevents attack amplification in MOON while maintaining nearly the same no-attack accuracy. Table 2 shows that we achieve higher post-attack accuracy across all heterogeneity levels. Our technique also outperforms FedAvg, except in a few scenarios. Techniques like MOON are designed to enhance accuracy under high heterogeneity ($\alpha = 0.1$). HYDRA-FL achieves a no-attack [attack] accuracy of 60.1[43.6], surpassing both MOON (58.8[39.9]) and FedAvg (57.76[40.9]).

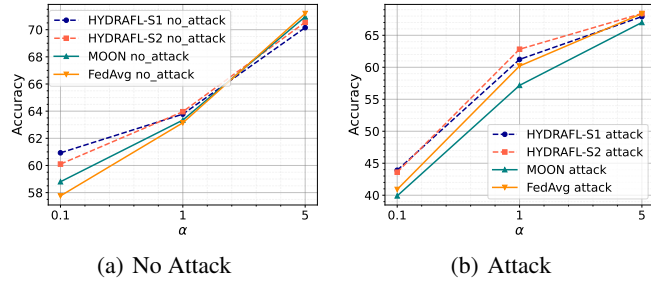


Figure 5: HYDRA-FL vs. MOON and FedAvg when auxiliary classifiers are placed at different shallow layers.

6 Analysis

In this section, we provide an in-depth analysis of HYDRA-FL. We begin with a qualitative analysis using t-distributed stochastic neighbor embedding (t-SNE [42]) plots to visualize the representations of the models. Then, we explore the impact of different design choices through ablation studies, focusing on the choice of the shallow layer for auxiliary classifiers and the distillation coefficients.

6.1 Qualitative Analysis

We show the t-SNE plots of the representations (Figure 6) generated by the client model for FedAvg, MOON, and HYDRA-FL for both attack and no-attack scenarios. The t-SNE plots show the classes as clusters. In the MOON attack scenario, the deviation from the no-attack scenario is much higher than the deviation between HYDRA-FL with and without attack, as evident from the spread of the class clusters, especially along the x-axis.

6.2 Ablation Study

Impact of choice of the shallow layer. Figure 5 illustrates the impact of the choice of the layer at which we insert our auxiliary classifier. We represent these choices by *HYDRAFL-S1* and *HYDRAFL-S2*, where the auxiliary classifier is inserted after the first and second convolutional layers, respectively.

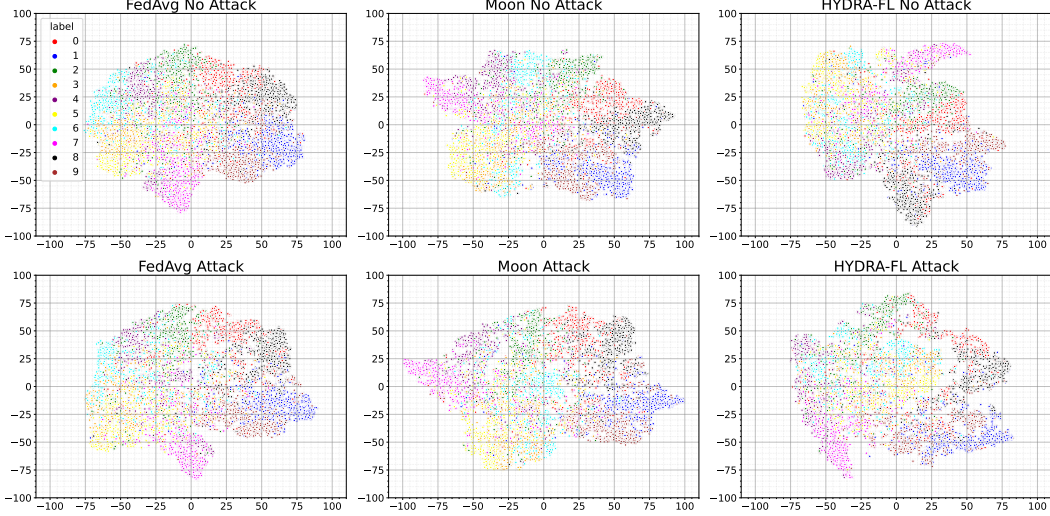


Figure 6: T-SNE visualizations of CIFAR10 on local model’s hidden representations ($\alpha = 0.5$) on FedAvg, MOON, and HYDRA-FL (ours). The attack vs. no-attack plot shows the deviation of the attack clusters from the no-attack clusters. Visually we can see MOON-attack has the greatest deviation, particularly along the x-axis, compared to FedAvg and HYDRA-FL.

We compare them in both attack and no-attack settings with simple MOON and FedAvg. In Figure 5(a), both HYDRAFL-S1 and HYDRAFL-S2 outperform other techniques at low heterogeneity in the absence of an attack but slightly underperform in low heterogeneity when $\beta = 5$. Figure 5(b) shows that both HYDRAFL-S1 and HYDRAFL-S2 achieve higher post-attack accuracy at all heterogeneity levels, with HYDRAFL-S2 giving a slightly higher accuracy than HYDRAFL-S1. The benefit from the contrastive loss reduces as we go shallower, so an optimal balance is necessary.

Impact of distillation coefficients. We examine the impact of distillation coefficients on the performance of FedNTD and HYDRA-FL. Figure 7 shows the post-attack accuracies with two different values of the *diminishing factor* $b = 1, 4$, resulting in output-layer NTD-loss coefficients of $\beta = 1$ and $\beta = 0.25$. Diminishing the coefficient β leads to improved performance, with a significant increase in post-attack accuracy for $\beta = 0.25$ at high heterogeneity ($\alpha = 0.05, 0.1$). As demonstrated in §3, β contributes to attack amplification in FedNTD. Reducing it while performing distillation at the auxiliary classifier yields the best performance. For example, at $\alpha = 0.05$, HYDRA-FL achieves 25.15% accuracy at $\beta = 1$, but a much higher accuracy of 28.81% at $\beta = 0.25$. Similar improvements are observed at other heterogeneity levels.

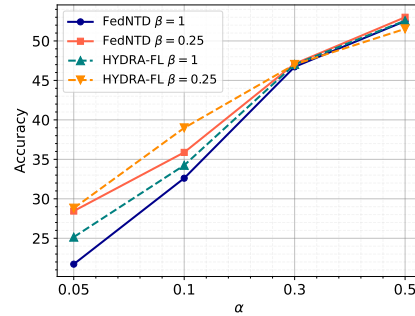


Figure 7: Comparison of performance of FedNTD-S with different values of β

7 Conclusion

In this paper, we first identified a critical issue in KD-based FL techniques that aim to tackle data heterogeneity: in the presence of model poisoning attacks, these techniques help the attacker amplify its effect, leading to reduced global model performance. We presented empirical evidence and theoretical reasoning to back this claim. This motivated us to propose HYDRA-FL: a hybrid knowledge distillation technique for robust and accurate FL technique that aims to tackle both data heterogeneity and model poisoning, two of the biggest problems in FL. Through extensive evaluation across three datasets and comparing with baseline techniques, FedNTD and MOON, we showed that HYDRA-FL achieves superior results.

References

- [1] How Apple personalizes Siri without hoovering up your data. <https://www.technologyreview.com/2019/12/11/131629/apple-ai-personalizes-siri-federated-learning/>.
- [2] Federated learning: Collaborative machine learning without centralized training data. <https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>, 2017.
- [3] Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How to backdoor federated learning. In *AISTATS*, 2020.
- [4] Moran Baruch, Baruch Gilad, and Yoav Goldberg. A Little Is Enough: Circumventing Defenses For Distributed Learning. In *NeurIPS*, 2019.
- [5] Arjun Nitin Bhagoji, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. Analyzing federated learning through an adversarial lens. In *ICML*, 2019.
- [6] Peva Blanchard, Rachid Guerraoui, Julien Stainer, et al. Machine learning with adversaries: Byzantine tolerant gradient descent. In *NeurIPS*, 2017.
- [7] X. Cao, J. Jia, Z. Zhang, and N. Gong. Fedrecover: Recovering from poisoning attacks in federated learning using historical information. In *2023 IEEE Symposium on Security and Privacy (SP) (SP)*, pages 326–343, Los Alamitos, CA, USA, may 2023. IEEE Computer Society.
- [8] Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. Provably Secure Federated Learning against Malicious Clients. In *AAAI*, 2021.
- [9] Hongyan Chang, Virat Shejwalkar, Reza Shokri, and Amir Houmansadr. Cronus: Robust and Heterogeneous Collaborative Learning with Black-Box Knowledge Transfer. *arXiv:1912.11279*, 2019.
- [10] Hong-You Chen and Wei-Lun Chao. Fedbe: Making bayesian model ensemble applicable to federated learning. *arXiv preprint arXiv:2009.01974*, 2020.
- [11] Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. Local Model Poisoning Attacks to Byzantine-Robust Federated Learning. In *USENIX*, 2020.
- [12] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [13] Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- [14] Peter Kairouz, H Brendan McMahan, Brendan Avent, et al. Advances and open problems in federated learning. *arXiv:1912.04977*, 2019.
- [15] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.
- [16] Momin Ahmad Khan, Virat Shejwalkar, Amir Houmansadr, and Fatima M Anwar. On the pitfalls of security evaluation of robust federated learning. In *2023 IEEE Security and Privacy Workshops (SPW)*, pages 57–68. IEEE, 2023.
- [17] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [18] Souvik Kundu, Qirui Sun, Yao Fu, Massoud Pedram, and Peter Beerel. Analyzing the confidentiality of undistillable teachers in knowledge distillation. *Advances in Neural Information Processing Systems*, 34:9181–9192, 2021.
- [19] Yann LeCun and Corinna Cortes. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.

- [20] Gihun Lee, Minchan Jeong, Yongjin Shin, Sangmin Bae, and Se-Young Yun. Preservation of the global knowledge by not-true distillation in federated learning. *Advances in Neural Information Processing Systems*, 35:38461–38474, 2022.
- [21] Lee-Gihun. Fedntd. <https://github.com/Lee-Gihun/FedNTD/tree/master>.
- [22] Daliang Li and Junpu Wang. Fedmd: Heterogenous federated learning via model distillation. *arXiv preprint arXiv:1910.03581*, 2019.
- [23] Liping Li, Wei Xu, Tianyi Chen, Georgios B Giannakis, and Qing Ling. RSA: Byzantine-robust stochastic aggregation methods for distributed learning from heterogeneous datasets. In *AAAI*, 2019.
- [24] Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. Federated learning on non-iid data silos: An experimental study. In *2022 IEEE 38th international conference on data engineering (ICDE)*, pages 965–978. IEEE, 2022.
- [25] Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10713–10722, 2021.
- [26] Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization. In *ICML*, 2021.
- [27] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2:429–450, 2020.
- [28] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019.
- [29] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. In *NeurIPS*, 2020.
- [30] Haoyu Ma, Tianlong Chen, Ting-Kuei Hu, Chenyu You, Xiaohui Xie, and Zhangyang Wang. Undistillable: Making a nasty teacher that cannot teach students. *arXiv preprint arXiv:2105.07381*, 2021.
- [31] Saeed Mahloujifar, Mohammad Mahmoody, and Ameer Mohammed. Universal multi-party poisoning attacks. In *ICML*, 2019.
- [32] H Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *AISTATS*, 2017.
- [33] El Mahdi El Mhamdi, Rachid Guerraoui, and Sébastien Rouault. The Hidden Vulnerability of Distributed Learning in Byzantium. In *ICML*, 2018.
- [34] Thomas Minka. Estimating a Dirichlet distribution, 2000.
- [35] Luis Muñoz-González, Battista Biggio, Ambra Demontis, Andrea Paudice, Vasin Wongrasamee, Emil C Lupu, and Fabio Roli. Towards poisoning of deep learning algorithms with back-gradient optimization. In *AISec*, 2017.
- [36] PyTorch Documentation. <https://pytorch.org/>, 2019.
- [37] QinbinLi. Moon. <https://github.com/QinbinLi/MOON/tree/main>.
- [38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [39] Sashank J Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and Hugh Brendan McMahan. Adaptive Federated Optimization. In *ICLR*, 2020.

- [40] V. Shejwalkar, A. Houmansadr, P. Kairouz, and D. Ramage. Back to the drawing board: A critical evaluation of poisoning attacks on production federated learning. In *2022 IEEE Symposium on Security and Privacy (SP)* (SP), pages 1117–1134, Los Alamitos, CA, USA, may 2022. IEEE Computer Society.
- [41] Virat Shejwalkar and Amir Houmansadr. Manipulating the Byzantine: Optimizing Model Poisoning Attacks and Defenses for Federated Learning. In *NDSS*, 2021.
- [42] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [43] Utilization of FATE in Risk Management of Credit in Small and Micro Enterprises. <https://www.fedai.org/cases/utilization-of-fate-in-risk-management-of-credit-in-small-and-micro-enterprises/>, 2019.
- [44] Cong Xie, Oluwasanmi Koyejo, and Indranil Gupta. Generalized byzantine-tolerant sgd. *arXiv:1802.10116*, 2018.
- [45] Cong Xie, Sanmi Koyejo, and Indranil Gupta. Fall of empires: Breaking Byzantine-tolerant SGD by inner product manipulation. *arXiv:1903.03936*, 2019.
- [46] Yueqi Xie, Weizhong Zhang, Renjie Pi, Fangzhao Wu, Qifeng Chen, Xing Xie, and Sunghun Kim. Robust federated learning against both data heterogeneity and poisoning attack via aggregation optimization. *arXiv preprint*, 2022.
- [47] Dong Yin, Yudong Chen, Kannan Ramchandran, and Peter Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. In *ICML*, 2018.
- [48] Lin Zhang, Li Shen, Liang Ding, Dacheng Tao, and Ling-Yu Duan. Fine-tuning global model via data-free knowledge distillation for non-iid federated learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10174–10183, 2022.
- [49] Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3713–3722, 2019.
- [50] Zaixi Zhang, Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. Fldetector: Defending federated learning against model poisoning attacks via detecting malicious clients. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 2545–2555, 2022.
- [51] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.
- [52] Yanlin Zhou, George Pu, Xiyao Ma, Xiaolin Li, and Dapeng Wu. Distilled one-shot federated learning. *arXiv preprint arXiv:2009.07999*, 2020.
- [53] Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. Data-free knowledge distillation for heterogeneous federated learning. In *International conference on machine learning*, pages 12878–12889. PMLR, 2021.

Appendix

We provide additional information for our paper, HYDRA-FL: Hybrid Knowledge Distillation for Robust and Accurate Federated Learning, in the following order:

- Limitations and Future Work (Appendix A)
- Terminology/Techniques (Appendix B)
- Adversarial Settings (Appendix C)
- Experimental Setup (Appendix D)
- Additional Results (Appendix E)

A Limitations and Future Work

Federated Learning can have very diverse setups, especially FL in an adversarial setting. We can have many setup combinations as we can choose between different aggregation rules, attacks, defenses, datasets, data modalities, data distribution types, data heterogeneity levels, number of clients, etc. Therefore, evaluating against all combinations of these settings is well beyond the scope of one paper. Hence, for this paper, we chose only a few combinations of FL settings and tried our best to show that the problem we identified using two representative FL techniques will also exist in similar techniques. Similarly, we laid out our solution as a general framework to achieve good performance under high heterogeneity and model poisoning simultaneously. To show generalizability, we tailored it to our two representative techniques, but it would be interesting to see how our solution adapts to and performs with other FL techniques in future works. Also, we have only used unimodal, i.e., image datasets for our evaluations. This was done to stay consistent with the implementations of the techniques chosen for our case study, FedNTD and MOON. However, the language modality is becoming popular now, and multimodal models such as CLIP [38] are being widely used as they achieve superior performance by combining both image and language modalities. We hope to incorporate language and multimodal models in our future works.

B Terminology/Techniques

B.1 FedNTD

FedNTD [20] is a KD-based technique that tackles the problem of data heterogeneity in FL. They first demonstrate that Data Heterogeneity causes local models to forget out-distribution knowledge, i.e., the data samples not part of the client’s local data. Therefore, to preserve the out-distribution knowledge, they introduce not-true distillation, which basically modifies the loss function for the client model’s local objective. FedNTD’s loss function is given by:

$$\mathcal{L} = \mathcal{L}_{CE}(y_c, y) + \frac{\beta}{b} \mathcal{L}_{KL}(\tilde{y}_c, \tilde{y}_s) \quad (11)$$

Here y is the target label, y_c is the client model’s output, \tilde{y}_s and \tilde{y}_c are the client model’s and the server model’s not-true logits, respectively.

B.2 MOON

MOON [25] also aims to solve the problem of data heterogeneity in FL. They do so by reducing the distance between the representation learned by the local model with that of the global model. MOON’s loss function is given by:

$$\mathcal{L} = \mathcal{L}_{CE}(y_c, y) + \frac{\mu}{b} \mathcal{L}_{con}(z_c, z_s) \quad (12)$$

Here y is the target label, y_c is the client’s output, z_c is the representation from the client’s final layer, z_s is the representation from the server’s final layer, and y_s is the server model’s output.

B.3 Shallow Layer and Shallow Distillation

Shallow layer. in a neural network refers to one of the early layers close to the input, as opposed to deeper layers that are closer to the output. In the context of a deep learning model, shallow layers generally capture low-level features, such as edges in images or simple patterns in data, while deeper layers capture more complex, abstract representations.

Shallow distillation. is a technique used in KD where the knowledge transfer happens at a shallow layer of the neural network rather than at the final output layer. In traditional KD, the student model tries to mimic the teacher model’s output at the final layer. In shallow distillation, an additional distillation loss is applied at one of the shallow layers of the student model. This helps the student model learn intermediate representations from the teacher, providing a more comprehensive learning experience. By aligning these intermediate representations, the student model gains a more robust understanding of the data, leading to better *generalization*.

Robustness against poisoning. Shallow layers are less affected by adversarial attacks that target the final output of the model. Applying distillation at a shallow layer reduces the impact of a poisoned global model because the knowledge transferred is more fundamental and less influenced by the adversarial manipulations that typically affect the deeper layers.

C Adversarial Settings

Here we present the details of the adversarial settings of our experiments. We explain our threat model, which attacks we are using and why we are using them, and the defense we are using.

C.1 Threat Model

Goal: Our untargeted poisoning adversary controls m out of N clients to manipulate the global model to misclassify all the inputs it can during testing. Unless stated otherwise, we assume 20% malicious clients. Most defense works assume high percentages of malicious clients to demonstrate that their defenses work even in highly adversarial settings. Hence, although unreasonable in practical FL settings [40], we follow prior defense works and use 20% malicious clients.

Knowledge: Following most of the defense works, we assume that the adversary knows the robust AGR that the server uses. As assumed by most works, the adversary knows the server’s AGR. To test the efficacy of our technique with a strong adversary, we consider the case where the adversary has access to not only the malicious clients’ data but also the benign clients’ data. This enables us to determine the upper bound of the efficacy of our technique.

Capabilities: Our adversary is strong enough to directly manipulate model updates of the malicious clients it controls. While poisoning attacks come in various types and flavors, we restrict ourselves to only model poisoning attacks. This is because model poisoning attacks are much stronger. It has been shown in [40] that model poisoning attacks are much stronger because they directly perturb the local model parameters. In contrast, data poisoning attacks perturb the data, subsequently perturbing the local and global models upon aggregation. Poisoning attacks can also be classified based on their error specificity. If the goal is to misclassify certain classes only, then it is a *targeted attack* and is often achieved by inserting a backdoor in the model that activates only for certain inputs. On the other hand, an *untargeted attack* indiscriminately lowers the accuracy for all inputs.

C.2 Attacks we use in our evaluation

We use two model poisoning attacks for our evaluations. By testing which attack worked well, we chose the Stat-Opt attack for MOON and the Dyn-Opt attack for FedNTD. Below, we briefly explain how they work:

- **Stat-Opt [11]:** gives an untargeted model poisoning framework and tailors it to specific defenses such as TrMean [47], Median [47], and Krum [6]. The adversary first calculates the mean of the benign updates, ∇^b , and finds the *static* malicious direction $w = -\text{sign}(\nabla^b)$. It directs the benign average along the calculated direction and scales it with γ to obtain the final poisoned update, $-\gamma w$.

- **Dyn-Opt [41]:** also gives an untargeted model poisoning framework and tailors it to specific defenses, similar to Stat-Opt but differs in the *dynamic* and *data-dependent* nature of the perturbation. The attack first computes the mean of benign updates, ∇^b , and a data-dependent direction, w . The final poisoned update is calculated as $\nabla' = \nabla^b + \gamma w$, where the attack finds the largest γ that can bypass the AGR. They compare their attack with Stat-Opt and show that the dataset-tailored w and optimization-based scaling factor γ make their attack much stronger.

C.3 Defense we use in our evaluation

We use the Trimmed Mean defense in our evaluations. Trimmed Mean [47, 44] is a foundational defense used in advanced AGRs [7, 50, 41]. The server receives model updates from each client, sorts each input dimension j , discards the m largest and smallest values (where m indicates malicious clients), and averages the rest.

D Experimental Setup

Models: For MOON, we use a base encoder with two 5×5 convolutional layers, each followed by a 2×2 max pooling layer and two fully connected layers with ReLU activation. The base encoder is followed by a projection head with an output dimension of 256. For FedNTD, we use a model (similar to the one in [32]) having two convolutional layers followed by a linear layer and a classification layer. For FedNTD, we test with different values and settle upon a diminishing factor $b = 1$ and $\gamma = 2$. For MOON, we set $\beta = 0$ and set $\gamma = 1$. We used PyTorch [36] for our implementation on an 8GB NVIDIA RTX 3060 Ti GPU. Each run of FedNTD and MOON took about 2-3 hours on our machine.

FL Settings: For FedNTD, we use 100 clients with a sampling ratio of 0.1, i.e., 10 clients are selected every round. We use momentum SGD with an initial learning rate of 0.1, weight decay of $1 \times e^{-5}$, batch size of 50, and momentum of 0.9. Each run consists of 200 rounds with 5 local epochs. For MOON, we use 10 clients with a sampling ratio of 1. We use SGD with an initial learning rate of 0.01, weight decay of $1 \times e^{-5}$, batch size of 64, and momentum of 0.9. Each run consists of 30 rounds with 10 local epochs, sufficient for convergence.

Data Partitioning: We use the widely used Dirichlet [34] distribution to generate the non-IID partitioning of data between clients. Dirichlet distribution works by sampling $p_k \sim \text{Dir}_N(\alpha)$ and assigns $p_{k,j}$ proportion of samples of class k to client j . A lower value of α corresponds to a higher level of heterogeneity since it means that most of the samples of a certain class belong to one client. Conversely, at a higher value of α , the class samples are more evenly distributed between the clients. Also, a characteristic of the Dirichlet distribution is that both local dataset size and local per-class distribution vary across clients.

Datasets: The three datasets we use in our experiments are:

- **MNIST [19]:** MNIST is a 10-class digit image classification dataset, which contains 70,000 grayscale images of size 28×28 . We divide all data among FL clients (100 for FedNTD and 10 for MOON) using the Dirichlet [39] distribution.
- **CIFAR10 [17]:** CIFAR10 is a 10-class classification task with 60,000 total RGB images, each of size 32×32 . Each class has 6000 training images and 1000 testing images. We divide all the data among 100 clients using the Dirichlet distribution, a popular synthetic strategy to generate FL datasets.
- **CIFAR100 [17]:** CIFAR100 is similar to CIFAR10, except that it is a 100-class classification task where each class has 600 images of size 32×32 . There are 500 training images and 100 test images per class. Like other datasets, we also partition this dataset using the Dirichlet distribution.

E Additional Results

In this section, we present some of the additional results we have obtained.

Table 3: FedNTD

Dataset	MNIST		CIFAR10								CIFAR100	
			0.05		0.1		0.3		0.5			
Techniques	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>	<i>no attack</i>	<i>attack</i>
Fedavg	92.12	74.48	44.69	31.27	54.67	35.67	66.34	42.53	70.57	48.27	26.17	12.92
MOON	93.03	58.09	46.94	21.72	56.95	32.61	68	46.72	71.79	52.51	29.1	13.92
Ours	92.69	76.67	46.92	25.15	57.12	34.25	68.1	47.03	71.22	52.57	28.9	14.33

E.1 FedNTD

For visual symmetry, we did not include the full table in §5, but we had also run our FedNTD experiments at $\alpha = 0.3$. We show the full FedNTD results in Table 3. Here, we can see that at $\alpha = 0.3$ too, we achieve superior results FedAvg and FedNTD in both benign and adversarial conditions.

E.2 MOON

We also ran ablation with MNIST for different shallow layers and diminishing coefficients. We show the results in Table 4, where we can see that at a lower μ , i.e., higher diminishing factor, we achieve the best results. A lower μ does give us better no-attack accuracy, but we lose a lot in the attack scenario.

Method	μ	no-attack	attack
HYDRA-FL s1	1	94.41	68.68
HYDRA-FL s2	1	91.78	68.13
HYDRA-FL s1	0.3	92.03	72.35
HYDRA-FL s2	0.3	92.92	73.55
HYDRA-FL s1	0.1	92.04	76.65
HYDRA-FL s2	0.1	93.93	72.54

Table 4: Comparison of HYDRA-FL for MOON with different distillation coefficients.

NeurIPS Paper Checklist

1. Claims

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

Answer: [\[Yes\]](#)

Justification: Yes, we have ensured that the main claims in the abstract and introduction accurately reflect the paper's contributions and scope.

Guidelines:

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

2. Limitations

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

Answer: [\[Yes\]](#) .

Justification: Yes, we have discussed the limitations and future work in Appendix A.

Guidelines:

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

3. Theory Assumptions and Proofs

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

Answer: [\[NA\]](#) .

Justification: There is no theoretical result in this paper that requires a full set of assumptions and correct proof.

Guidelines:

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

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Yes, we have fully disclosed the information needed to reproduce the main experimental results of the paper. They are written in Section 5 and Appendix D.

Guidelines:

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

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes] .

Justification: Yes, we are submitting the code for HYDRA-FL in the supplementary material. For now, we have given the code where HYDRA-FL is adapted to FedNTD. We provide an "instructions.txt" file to reproduce our results. We will publish our full code for FedNTD and MOON on github with the final version of this paper.

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, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes] .

Justification: We specify the training and test details in Appendix D.

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: [No] .

Justification: We did not have enough compute resources to completely re-run all the experiments for different seeds and report error bars for different runs. We are currently re-running the error bar experiments, and we plan to include all the experiments with different seeds in the final version.

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.
- 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.

8. Experiments Compute Resources

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

Answer: [Yes] .

Justification: We present these details in Appendix D.

Guidelines:

- 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).

9. Code Of Ethics

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

Answer: [Yes] .

Justification: Yes, to the best of our knowledge, our paper conforms to the NeurIPS Code of Ethics in every aspect.

Guidelines:

- 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).

10. Broader Impacts

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

Answer: [NA] .

Justification: Or work does not have such a societal impact that requires discussion in the paper.

Guidelines:

- 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 generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

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

Answer: [NA] .

Justification: To the best of our knowledge, our paper poses no such risks. We use publicly available code and data for our work.

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes] .

Justification: We have cited all three datasets; MNIST [19], CIFAR10 [17], and CIFAR100 [17]. Their licenses are not mentioned on paperswithcode.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

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

Answer: [Yes] .

Justification: Yes, we are submitting the code for HYDRA-FL in the supplementary material. For now, we have given the code where HYDRA-FL is adapted to FedNTD. We provide an "instructions.txt" file to reproduce our results.

Guidelines:

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

14. Crowdsourcing and Research with Human Subjects

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

Answer: [NA] .

Justification: Our paper does not involve any crowdsourcing experiments nor research with human subjects.

Guidelines:

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

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

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

Answer: [NA] .

Justification: Our paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

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