

---

# UniFed: A Benchmark for Federated Learning Frameworks

---

<b>Xiaoyuan Liu</b> UC Berkeley xiaoyuanliu@berkeley.edu	<b>Tianneng Shi</b> Shanghai Jiao Tong University stneng@sjtu.edu.cn	<b>Chulin Xie</b> UIUC chulinx2@illinois.edu
<b>Qinbin Li</b> National University of Singapore qinbin@comp.nus.edu.sg	<b>Kangping Hu</b> Zhejiang University hukp@zju.edu.cn	<b>Haoyu Kim</b> Peking University jhy2001@pku.edu.cn
<b>Xiaojun Xu</b> UIUC xiaojun3@illinois.edu	<b>Bo Li</b> UIUC lbo@illinois.edu	<b>Dawn Song</b> UC Berkeley dawnsong@cs.berkeley.edu

## Abstract

1 Federated Learning (FL) has become a practical and popular paradigm in machine  
2 learning. However, currently, there is no systematic solution that covers diverse  
3 use cases. Practitioners often face the challenge of how to select a matching FL  
4 framework for their use case. In this work, we present UniFed, the first unified  
5 benchmark for standardized evaluation of the existing open-source FL frameworks.  
6 With 15 evaluation scenarios, we present both qualitative and quantitative eval-  
7 uation results of nine existing popular open-sourced FL frameworks, from the  
8 perspectives of functionality, usability, and system performance. We also provide  
9 suggestions on framework selection based on the benchmark conclusions and point  
10 out future improvement directions.

## 11 1 Introduction

12 Federated Learning (FL) [42, 26] has become a practical and popular paradigm for training machine  
13 learning (ML) models. There are many existing open-source FL frameworks. However, unlike  
14 Pytorch [43] and TensorFlow [8] for ML, currently, there is not a dominant systematic solution that is  
15 maturely developed for most use cases.

16 To compare existing open-source solutions, we created UniFed, an FL benchmark for standardized  
17 evaluations of FL frameworks. Specifically, UniFed helps answer two questions:

- 18 • How to qualitatively and quantitatively characterize an FL framework?
- 19 • How to choose the best FL framework for a specific real-world application?

20 We find that existing FL frameworks have significant qualitative differences which we present in Table  
21 2 and Table 3. In addition, our training experiments on nine existing FL frameworks with different  
22 FL algorithm implementations suggest that the selection of model type is the main factor that affects  
23 model performance compared with the selection of algorithm and framework. Our measurement of  
24 system performance shows that, interestingly, when considering training efficiency, communication  
25 efficiency, and memory usage, there is no framework that consistently outperforms others.

26 The contribution of this paper is summarized below.

- 27 1. We define the criteria to characterize an FL framework, including functionality support,  
28 system performance, and usability. We also develop a toolkit where users can easily deploy  
29 and test FL frameworks in various settings in one command, which facilitates a workflow of  
30 a standardized quantitative evaluation.
- 31 2. We collect and categorize a list of nine representative FL frameworks. With the criteria and  
32 the toolkit, we benchmark and compare them both qualitatively and quantitatively.
- 33 3. Summarizing the result in our evaluation, we present a complete guideline that helps FL  
34 practitioners choose the FL framework for a specific real-world application.

## 35 2 Related work

### 36 2.1 Existing datasets and frameworks

37 We give background on existing FL datasets in Appendix [A](#). In our benchmark, we cover both the  
38 simulated datasets and real federated datasets from diverse application domains for evaluation. Specif-  
39 ically, considering the real-world usage of FL frameworks, we adopt the datasets from the LEAF [\[10\]](#)  
40 for experiments on cross-device horizontal FL, which are more practical than the simulated datasets.  
41 For cross-silo horizontal FL and vertical FL, we adopt the generated datasets from FATE [\[38\]](#) for  
42 evaluation, which cover representative FL applications among institutions from finance to healthcare.

43 There are many system construction efforts on building frameworks to support various FL scenarios.  
44 In this work, we focus on open-source FL frameworks that are available for evaluation. Here we  
45 identify three general categories along with representative examples.

46 **All-in-one frameworks.** Great efforts have been made to construct a single framework that covers  
47 most FL-related techniques in both horizontal and vertical FL settings. Such FL frameworks (FATE  
48 [\[38\]](#), FedML [\[22\]](#), PaddleFL [\[6\]](#), Fedlearner [\[3\]](#)) focus on the coverage of the functionalities and are  
49 often constructed with great engineering efforts. For example, FATE from WeBank is an industrial-  
50 grade FL framework that aims to provide FL services for enterprises and institutions.

51 **Horizontal-only frameworks.** Instead of aiming to support diverse applications with both horizontal  
52 and vertical FL, some frameworks (TFF [\[7\]](#), Flower [\[9\]](#), FLUTE [\[16\]](#)) aim to provide easy-to-use  
53 APIs for users to adopt and develop horizontal FL algorithms. For example, based on TensorFlow [\[8\]](#),  
54 TFF [\[7\]](#) from Google provides federated learning API and federated core API for users to apply and  
55 design FL algorithms, respectively.

56 **Specialized frameworks.** While the above frameworks support the general development of FL, some  
57 frameworks (CrypTen [\[27\]](#), FedTree [\[34\]](#)) are specially designed for specific purposes. CrypTen [\[27\]](#)  
58 focuses on providing secure multi-party computation [\[46\]](#) primitives, while FedTree is designed for  
59 the federated training of decision trees.

### 60 2.2 Existing benchmarks

61 We give background on existing FL benchmarks in Appendix [B](#). In summary, they mainly focus  
62 on creating federated datasets in different tasks, either from natural client data or from artificially  
63 partitioned centralized datasets, to evaluate FL *algorithms*. However, they do not provide the  
64 systematic evaluation of FL *frameworks* that are built with industry efforts and used as real FL  
65 systems in practice. To fill in this gap, we benchmark nine open-source FL frameworks with 15  
66 common FL datasets to cover different FL settings, data modalities, task, as well as workload sizes.

## 67 3 Benchmark design

68 To get first-hand experience and quantitatively evaluate target FL frameworks, we design UniFed  
69 benchmark toolkit and integrate all target frameworks with minimal intrusions in coding. Figure [1](#)  
70 shows an example evaluation workflow using our toolkit. UniFed toolkit contains four components:

71 (1) An environment launcher that provides a Command Line Interface (CLI) to read the experiment  
 72 specification from a configuration file and launch a distributed testing environment. (2) A scenario  
 73 loader python package that facilitates easy access to evaluation scenarios with automatic caching. (3)  
 74 A set of framework-dependent code patches that inject extra code to the target FL frameworks for  
 75 evaluation-related data loading and performance logging. (4) A global log analyzer that collects log  
 76 files from distributed evaluation nodes and reports the benchmarking result. With our benchmark  
 77 toolkit, one can start an FL training experiment using any FL framework with a one-line command.

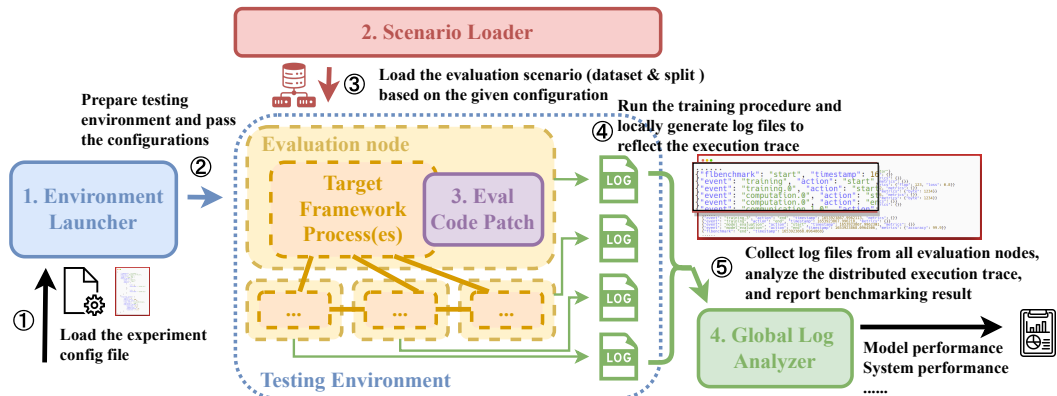


Figure 1: Design of the benchmark workflow.

78 To incorporate a new FL framework into our workflow, we first update our environment launcher to  
 79 support the deployment of the target framework based on the instructions from its documentation.  
 80 We then write code patches to the target framework for integration: (a) Support the data loading  
 81 of our evaluation scenarios. (b) Configure the framework behavior based on the configuration file  
 82 received from the launcher. (c) Generate timestamped event logs to record the training procedure  
 83 so that different frameworks generate log files with a unified format for fair comparisons. More  
 84 details are given in Section 5.1. To evaluate the target framework, following the steps in Figure 1,  
 85 one needs to write a configuration file specifying the experiment details (e.g. which framework to  
 86 use, what model type to train, what learning rate to apply) and start the training with the CLI from  
 87 the environment launcher. One can also run the CLI provided by the global log analyzer anytime  
 88 to see if the training has finished and, after the training finishes, get a comprehensive report of both  
 89 model and system performance.

90 Overall, users can easily adopt our UniFed toolkit to deploy and test different FL frameworks  
 91 in comprehensive scenarios with a one-line command. Moreover, developers can integrate their  
 92 frameworks into our toolkit for comparison following the same workflow.

## 93 4 Evaluation principles of UniFed

### 94 4.1 Evaluation scenarios and evaluation targets

95 For clarity, we use the term model performance to refer to the model’s ability to perform a task  
 96 and the term system performance to refer to its training efficiency (mainly discussed in Section  
 97 4.3). Theoretically, despite the difference in implementation (e.g. different ML backends and  
 98 communication orchestration), as long as the frameworks achieve the same mathematical procedure  
 99 with the same FL training algorithm, the resulting model should have similar model performance. In  
 100 this work, we measure the model performance across different frameworks to verify this statement.  
 101 We also explore the algorithm efficiency difference across different training algorithms and models.

102 We first define the scope of an evaluation scenario in our benchmark. Inspired by experiments in  
 103 existing FL studies, in this work, we define *one evaluation scenario* to be a set of clients who hold  
 104 fixed dataset splits (training/validation/test) with a fixed partition across different clients and an  
 105 optional aggregator/arbiter who potentially also holds a dataset with only the test set. In vertical FL,  
 106 each data instance has a unique identifier (id) and vertical split data instances have aligned identifiers

Setting	Scenario name	Modality	Task type	Performance metrics	Client number	Sample number
cross-device horizontal	celeba [39]	Image	Binary Classification (Smiling vs. Not smiling)	Accuracy	894	20,028
	femnist [30,14,10]	Image	Multiclass Classification (62 classes)	Accuracy	178	40,203
	reddit [10]	Text	Next-word Prediction	Accuracy	813	27,738
cross-silo horizontal	breast_horizontal [1]	Medical	Binary Classification	AUC	2	569
	default_credit_horizontal [47,17]	Tabular	Binary Classification	AUC	2	22,000
	give_credit_horizontal [4]	Tabular	Binary Classification	AUC	2	150,000
	student_horizontal [15,17]	Tabular	Regression (Grade Estimation)	MSE	2	395
	vehicle_scale_horizontal [44,17]	Image	Multiclass Classification (4 classes)	Accuracy	2	846
<b>Vertical split details</b>						
vertical	breast_vertical [1]	Medical	Binary Classification	AUC	A: 10 features 1 label B: 20 features	
	default_credit_vertical [47,17]	Tabular	Binary Classification	AUC	A: 13 features 1 label B: 10 features	
	dvisits_vertical [11]	Tabular	Regression (Number of consultations Estimation)	MSE	A: 3 features 1 label B: 9 features	
	give_credit_vertical [4]	Tabular	Binary Classification	AUC	A: 5 features 1 label B: 5 features	
	motor_vertical [2]	Sensor data	Regression (Temperature Estimation)	MSE	A: 4 features 1 label B: 7 features	
	student_vertical [15,17]	Tabular	Regression (Grade Estimation)	MSE	A: 6 features 1 label B: 7 features	
	vehicle_scale_vertical [44,17]	Image	Multiclass Classification (4 classes)	Accuracy	A: 9 features 1 label B: 9 features	

Table 1: Evaluation scenarios in UniFed. UniFed borrow 15 datasets from existing works to cover different FL settings, modalities, task types, and workload sizes.

107 in different participants. The model performance is measured with all available test instances among  
108 participants in an aligned, unweighted, but non-deduplicated way. We list all evaluation scenarios  
109 considered in this paper in Table 1.

110 We then define the scope of an evaluation target that determines the granularity of our evaluation. FL  
111 frameworks often support multiple FL algorithms (e.g. FedAvg [42], SecureBoost [13]) that uses  
112 different local training methods and different mathematical procedures for aggregation. Using a  
113 specific FL algorithm, one can train different ML models (e.g. linear regression, logistic regression,  
114 MLP, LeNet [31]) for a given task. Moreover, with different training parameters (e.g. epoch number,  
115 batch size, learning rate, choice optimizer), the same ML model can achieve different final model  
116 performances for the same task. Considering all differences, in this work, we define the basic unit for  
117 evaluation in the benchmark to be a combination of (*FL framework, FL algorithm, ML model*) and  
118 always measure the performance with a proper set of hyperparameters, which are chosen separately  
119 for different evaluation scenarios in advance with grid searches for the best model performance.

## 120 4.2 Functionality support

121 Different FL frameworks support different sets of functionalities. Specifically, we consider the model  
122 support in both horizontal and vertical settings, the deployment support, and privacy-protection  
123 features. The model support reflects whether the evaluation target is able to train a specific type of  
124 model in a specific setting. The deployment support measures the scalability of the evaluation target  
125 and challenges its communication infrastructure. And for privacy-protection features, we examine  
126 whether the evaluation target has proper mechanisms to resist certain types of privacy threats. We  
127 give a functionality comparison for all nine frameworks in Table 2. Note that, UniFed focus on  
128 features that are commonly supported by the frameworks off the shelf. There are also latest research  
129 projects developing more functionalities for better training optimization, better robustness [50], more  
130 comprehensive differential privacy, and improved fairness. Most frameworks in our evaluation can  
131 potentially be extended for those additional functionalities, which is an interesting future direction.

132 From the table, we make the following observations.

Framework	All-in-one frameworks				Horizontal-only frameworks			Specialized frameworks	
	FATE	FedML	PaddleFL	Fedlearner	TFF	Flower	FLUTE	CrypTen	FedTree
<b>Model support - Horizontal</b>									
Regression	Y	Y	Y	Y	Y	Y	Y	N/A	N
Neural network	Y	Y	Y	Y	Y	Y	Y	N/A	N
Tree-based model	Y	N	N	N	N	N	N	N/A	Y
<b>Model support - Vertical</b>									
Regression	Y	Y	Y*	N	N	N	N	Y	N
Neural network	Y	N	Y*	Y	N	N	N	Y	N
Tree-based model	Y	N	N	Y	N	N	N	N	Y
<b>Deployment support</b>									
Single-host simulation	Y	Y	Y	Y	Y	Y	Y	Y	Y
Multi-host deployment (<16 hosts)	Y	Y*	Y	Y	N	Y	Y	Y	Y
Cross-device deployment (>100 host)	N	Y*	Y	Y	N	Y	Y	N/A	Y
<b>Privacy protection against the semi-honest server</b>									
Does not require a 3rd party aggregator (vertical)	Y	N	Y	Y	N	N	N	Y	Y
Aggregator does not learn model param (arbitar scenario)	Y	N	Y	N	N	N	N	N/A	Y
Aggregator does not learn individual model gradient (secagg)	Y	Y	Y	N	Y	N	N	N/A	Y
<b>Privacy protection against semi-honest peer clients</b>									
Clients does not learn anything about the model param (vertical)	Y	N	Y	N	N/A	N/A	N/A	Y	Y
Clients does not learn gradients from other clients (vertical)	Y	Y	Y	N	N/A	N/A	N/A	Y	Y
<b>Privacy protection in the final model</b>									
Support training with central DP (dpsgd/gradient edits)	N	N	Y	N	Y	Y	Y	N	Y

Table 2: Functionality support in different FL frameworks. Asterisks indicate a claimed support for certain functionalities that are missing or cannot run in the open-source implementation.

133 1. **Model support.** For horizontal settings, most frameworks support both regression and neural  
134 networks, while only a few (FATE, FedTree) support tree-based models. For vertical settings, only all-  
135 in-one frameworks support the corresponding algorithms and the coverage is incomplete. Tree-based  
136 vertical training is only supported by three frameworks (FATE, Fedlearner, and FedTree).

137 2. **Deployment support.** While all frameworks support the single-host deployment as a basic  
138 functionality, surprisingly most frameworks provide the multi-host deployment option for realistic  
139 FL simulation. The only exception is TensorFlow Federated which has multi-host deployment as  
140 its incoming feature in development. However, for cross-device support where we challenge the  
141 scalability of the evaluation target, although most frameworks claim they support the cross-device  
142 training, we experience various glitches in practice that prevent a successful deployment. More  
143 details are discussed in Section 5.

144 3. **Privacy enhancement.** We investigate the privacy features that are actually implemented in our  
145 evaluation targets and categorize them based on their different threat models. Aligned with [26], we  
146 identify three types of protections against attackers with different access. (1) Specifically, to keep  
147 private information from an honest-but-curious central server in a vertical setting, some frameworks  
148 (FATE, PaddleFL, Fedlearner, CrypTen, FedTree) support different protocols without arbiters which  
149 provide the ultimate protection. For example, FATE uses HE-based solutions [20, 48, 49] for  
150 regression and neural network while CrypTen uses sMPC-based solutions [28]. For tree-based  
151 models, most frameworks use SecureBoost [13] in the vertical setting and, in the horizontal setting, a  
152 histogram secure aggregation (HistSecAgg) mentioned in [5]. In addition, some frameworks (FATE,  
153 PaddleFL, FedTree) take advantage of arbiters for better computation efficiency but do not reveal any  
154 model parameter. In settings where the aggregator needs the final model as the output, there is also the  
155 option of secure aggregation that prevents the aggregator from learning individual model gradients.  
156 (2) On the other side, to prevent clients from getting extra information in vertical settings, most  
157 frameworks that support vertical settings implement corresponding protection. The only exception is  
158 Fedlearner which only implements split learning and introduces certain amounts of gradient leakage  
159 [33]. We also notice that most implemented protection mechanisms are assuming a semi-honest  
160 model. (3) Finally, to protect user privacy and defend potential privacy attacks (e.g., membership  
161 inference, model-inversion) on the final production model, some frameworks (PaddleFL, TFF, Flower,  
162 FLUTE, FedTree) support applying differential privacy [18] in training.

### 163 4.3 System performance

164 Although there is a huge overlap in the functionality support for different FL frameworks, the  
165 implementations are often quite different, leading to different performance characteristics. To finish  
166 the complete FL training task, the frameworks often need to preprocess the data, locally computes  
167 certain functions, and potentially communicate with an aggregator to collaboratively learn the model.  
168 In most cases, the frameworks improve the model iteratively and repeat the above steps after a  
169 configured number of epochs or until certain criteria are matched.

Framework	All-in-one frameworks				Horizontal-only frameworks			Specialized frameworks	
	FATE	FedML	PaddleFL	Fedlearner	TFF	Flower	FLUTE	CrypTen	FedTree
<b>Documentation</b>									
Detailed tutorial	Y	Y	Y	N	Y	Y	Y	Y	Y
Code example	Y	Y	Y	Y	Y	Y	Y	Y	Y
API documentation	Y	N	N	N	Y	Y	N	Y	Y
<b>Engineering</b>									
Native test & benchmark	Y	Y	N	N	Y	Y	Y	Y	Y
GPU support	Y	Y	Y	Y	Y	Y	Y	Y	Y
<b>Built-in ML building block</b>									
CNN	Y	Y	Y	Y	Y	Y	Y	Y	N
RNN	Y	Y	Y	Y	Y	Y	Y	N	N
Rich Optimizers	Customized	Torch	PaddlePaddle	TensorFlow	TensorFlow	Y	Torch	Only SGD	N/A

Table 3: Usability feature comparison in different FL frameworks.

170 In this work, we target to measure the system performance in three aspects: training efficiency,  
171 communication cost, and resource consumption. Specifically, we are interested in a direct comparison  
172 between different frameworks training the same ML model. Because of the difference in their  
173 ML backend, communication orchestration, and implementation quality for model aggregation and  
174 synchronization, we target to find the best FL frameworks for each of our evaluation scenarios . We  
175 use logging in UniFed toolkit for performance tracking which we discuss in Section 3 to record and  
176 analyze the complete training procedure of a potentially distributed evaluation target. We discuss the  
177 evaluation result in Section 5 and a direct system performance comparison is given in Table 7.

#### 178 4.4 Usability

179 In addition to functionality and efficiency, whether the framework is easy to learn and convenient to  
180 use also affects its popularity. In this work, we first define a set of qualitative attributes and apply  
181 them to different FL frameworks to measure their usability. Specifically, we focus on three aspects:  
182 documentation, engineering, and built-in ML components. As frameworks often use their own term  
183 to refer to different pieces in their documentation, we standardize our requirement of tutorial, code  
184 example, and API documentation in Appendix C. In terms of engineering efforts, we mainly check if  
185 the target framework has its own tests and performance benchmark, and also check whether the GPU  
186 support can be explicitly configured. Last, we examine whether the framework has integrated basic  
187 ML building blocks like specific network structures and optimizers for convenient usage.

188 Based on the criteria listed above, we show our evaluation result in Table 3. Most frameworks  
189 provide details on their installation and usage. FedML, Fedlearner, and FLUTE do not provide  
190 API documentation for users to set up different FL scenarios easily. In terms of engineering efforts,  
191 all frameworks provide internal testing and benchmarking code except PaddleFL and Fedlearner.  
192 Moreover, all frameworks support the usage of GPU to accelerate training. In terms of ML building  
193 blocks, most frameworks have integrated CNN and RNN except CrypTen and FedTree, which are  
194 designed for specific purposes. FATE and Flower are compatible with different backend libraries  
195 such as TensorFlow and PyTorch, while the other frameworks support optimizers in its own backend.

## 196 5 Benchmark evaluation

### 197 5.1 Implementation

198 We implement and open source 1 UniFed toolkit discussed in Section 3. Specifically, our environment  
199 launcher uses SSH to connect to the evaluation node and prepares the testing environment. We wrap  
200 the data loading for datasets from [38], further automate the file caching from [10], and fix the dataset  
201 splits for the evaluation scenarios as discussed in Section 4.1. Our own logging format based on the  
202 JSON file structure records the timestamp for critical events. We explain the details of the logging  
203 format in Appendix D. For each framework, we create a separate code patch following the principle  
204 of minimal intrusion and resource consumption. We explain the details about separate patches to  
205 each individual framework in Appendix E. All experiments use evaluation nodes with 20 vCPU in  
206 Intel Xeon Gold 6230.

<https://github.com/AI-secure/FLBenchmark-toolkit>

Setting	Model	FATE	FedML	PaddleFL	Fedlearner	TFF	Flower	FLUTE
femnist cross-device (Accuracy)	logistic_regression	/	0.083	0.053	/	0.058	0.036	0.072
	mlp_128	/	0.652	0.591	/	0.644	0.663	0.641
	mlp_128_128_128	/	0.701	0.671	/	0.722	0.707	0.697
	lenet	/	0.822	0.792	/	0.822	0.819	0.820
give credit cross-silo (AUC)	logistic_regression	0.693	0.788	0.788	0.790	0.790	0.795	0.790
	mlp_128	0.830	0.832	0.828	0.834	0.832	0.831	0.833
	mlp_128_128_128	0.831	0.834	0.827	0.835	0.832	0.834	0.834

Table 4: FedAvg with different models on different tasks in the horizontal FL setting. FATE and Fedlearner do not support cross-device setting and are excluded from the comparison. We can observe that different FL frameworks show similar performance in general when using the same model.

## 207 5.2 Benchmark results

208 With UniFed toolkit, we run experiments and present representative benchmark quantitative results  
 209 related to the research questions in Section 1 and Section 4. We also analyze possible reasons for the  
 210 experiment outcomes by comparing implementations and designs in the frameworks.

211 **RQ1: Does the choice of FL framework affects the model performance trained using the**  
 212 **same FL algorithm?** As mentioned in Section 4.1, we expect a unified model performance across  
 213 different FL frameworks because of the same mathematical procedure of training. With our benchmark  
 214 result, we verify this finding with the most commonly supported FedAvg and the result is shown in  
 215 Table 4. Results for FATE and Fedlearner are partially missing due to their limited support in the  
 216 cross-device setting. Specifically, a file naming issue in FATE prevents it from scaling hundreds of  
 217 clients, and although Fedlearner supports cross-device training, it does not support sampling a subset  
 218 of clients and has different synchronizing mechanisms which prevent a fair comparison.

219 In Table 4, the model performance of the same model is generally consistent across different FL  
 220 framework implementations (performance difference within 1.1% in most cases) and the trend that  
 221 deeper models perform better can be observed in all frameworks for the selected scenarios. In addition,  
 222 we observe that (1) The logistic regression model does not work well in femnist scenario, which  
 223 leads to consistently poor performance in all frameworks. (2) In the cross-silo setting, the logistic  
 224 regression model in FATE has a relatively low performance which might be relevant to its default  
 225 early-stop behavior triggered by convergence. (3) The PaddleFL model performance is unstable and  
 226 consistently lower, which is potentially caused by its different ML backend PaddlePaddle [41].

227 **RQ2: Are different FL algorithm implementations comparable when training the same type**  
 228 **of model?** In addition to FedAvg, FL frameworks also implemented other FL algorithms to cover  
 229 specific use cases. As the selection of algorithms is less consistent for tree-based and vertical cases,  
 230 in this research question, we focus on a comparison between different frameworks with different FL  
 231 algorithm implementations training the same model. The result is presented in Table 5. We note that  
 232 FedML only supports regression in the vertical setting. For PaddleFL, we failed to run the sMPC  
 233 example following the official instructions in its latest version 1.2.0 and its split learning support is  
 234 also removed. Fedlearner only provides one-layer networks for split learning off-the-shelf. CrypTen  
 235 does not support tree-based models and FedTree does not support non-tree-based models.

236 In Table 5, again we observe relatively consistent model performance in each row, which suggests the  
 237 model selection is still the main factor that influences the model performance even with different FL  
 238 algorithm implementations. In addition, to explain the larger diversity compared with Table 4, we note  
 239 that (1) For the logistic regression, FATE achieves slightly better model performance probably due to  
 240 its default regularization option, while FedML suffers a performance loss that might be caused by its  
 241 default LeakyReLU activation. (2) FL algorithm in FATE failed to efficiently support the training of a  
 242 3-layer multi-level perceptron (MLP). We report the result after an insufficient training of one epoch  
 243 which takes more than 8 hours. CrypTen achieves better performance with sufficient training using  
 244 an efficient sMPC-based approach. (3) Tree-based models have more consistent model performance  
 245 despite their different FL algorithm implementations in different programming languages.

Setting	Model	FATE	FedML	PaddleFL	Fedlearner	CrypTen	FedTree
default credit vertical (AUC)	Regression (logistic_regression)	0.717 HE-based	0.650 HE-based	/	/	0.708 sMPC-based	/
	Neural network (mlp_128_128_128)	0.737 (1 epoch) HE-based	/	/	/	0.789 sMPC-based	/
	Tree-based model (gbdt_64_64_6)	0.820 SecureBoost	/	/	0.819 SecureBoost	/	0.817 SecureBoost
give credit horizontal (AUC)	Tree-based model (gbdt_64_64_6)	0.861 HistSecAgg	/	/	/	/	0.861 HistSecAgg

Table 5: Comparison among different FL algorithm implementations that train the same model. We observe that the model performance is still mainly determined by the model selection.

Setting	Name	1st		2nd		3rd	
		alg&model	perf	alg&model	perf	alg&model	perf
cross-device horizontal	celeba (Accuracy)	FedAvg leaf_cnn	90.19%	FedAvg resnet_18	88.99%		
	feminist (Accuracy)	FedAvg lenet	82.23%	FedAvg mlp_128_128_128	72.24%	FedAvg mlp_128	66.33%
	reddit (Accuracy)	FedAvg lstm	13.36%				
cross-silo horizontal	breast_horizontal (AUC)	FedAvg mlp_128_128_128	98.86%	FedAvg logistic_regression	98.70%	FedAvg mlp_128	98.54%
	default_credit_horizontal (AUC)	HistSecAgg gbdt_64_64_6	78.46%	FedAvg mlp_128_128_128	77.70%	FedAvg mlp_128	77.21%
	give_credit_horizontal (AUC)	HistSecAgg gbdt_64_64_6	86.10%	FedAvg mlp_128_128_128	83.45%	FedAvg mlp_128	83.38%
	student_horizontal (MSE)	FedAvg mlp_128_128_128	21.04	FedAvg mlp_128	21.99	HistSecAgg gbdt_64_64_6	22.79
	vehicle_scale_horizontal (Accuracy)	FedAvg mlp_128_128_128	100.0%	FedAvg mlp_128	100.0%	HistSecAgg gbdt_64_64_6	99.64%
vertical	breast_vertical (AUC)	SecureBoost gbdt_64_64_6	100.0%	sMPC-based mlp_128_128_128	100.0%	sMPC-based mlp_128	99.97%
	default_credit_vertical (AUC)	SecureBoost gbdt_64_64_6	81.99%	sMPC-based mlp_128_128_128	78.89%	sMPC-based mlp_128	77.87%
	dvisits_vertical (MSE)	SecureBoost gbdt_64_64_6	0.32	sMPC-based mlp_128_128_128	0.57	sMPC-based mlp_128	0.60
	give_credit_vertical (AUC)	SecureBoost gbdt_64_64_6	86.79%	sMPC-based mlp_128_128_128	83.38%	sMPC-based mlp_128	82.79%
	motor_vertical (MSE)	sMPC-based mlp_128_128_128	3.66E-4	SecureBoost gbdt_64_64_6	3.64E-3	sMPC-based mlp_128	9.98E-03
	student_vertical (MSE)	SecureBoost gbdt_64_64_6	3.26	sMPC-based mlp_128_128_128	11.03	sMPC-based mlp_128	12.43
	vehicle_scale_vertical (Accuracy)	SecureBoost gbdt_64_64_6	99.17%	sMPC-based mlp_128_128_128	96.34%	sMPC-based mlp_128	94.21%

Table 6: Best algorithm and model combinations for each evaluation scenario. Tree-based models generally have advantages in vertical settings and deeper models are often preferred.

246 **RQ3: How to select a model and FL algorithm combination to achieve a good model perfor-**  
247 **mance for the given application scenario?** From RQ1 and RQ2, we verified that the model type  
248 is the major factor that influences the model performance as long as the implementation in the FL  
249 framework is correct. In RQ3, we want to find the best of such combination among all available  
250 FL algorithms and model combinations we tested in UniFed evaluation scenarios. Specifically, for  
251 cross-silo horizontal and vertical settings, we compare available models of regression, shallow neural  
252 network (1-layer MLP), deep neural network (3-layer MLP), and tree-based model (GBDT). For  
253 the cross-device settings, we find promising models that are available off-the-shelf (CNN model in  
254 LEAF [10], ResNet [23], LeNet [32], LSTM [24]) for a reference. We present a ranked comparison  
255 result in Table 6.

256 We notice that in some scenarios, the model performance is sensitive to the change in model type,  
257 while for other scenarios, the difference is less significant and sometimes multiple models get good  
258 performance. Specifically, (1) Tree-based models often perform better in vertical settings, in some  
259 cases even by a large margin. (2) Deeper neural networks often achieve better performance than  
260 shallow ones in most cases. We recommend the practitioners find scenarios in the benchmark that  
261 are similar to their use case for a reference for the model selection. For scenarios where the model  
262 performance is less sensitive to the model selection, the practitioners should consider comparing the  
263 system performance, which is discussed in the next RQ.



Setting	Model	FATE	FedML	PaddleFL	Fedlearner	TFF	Flower	FLUTE	CrypTen	FedTree
<b>Training time in total</b>										
femnist horizontal	Neural network (lenet)	/	2,000 epochs 2,373 s FedAvg	2,000 epochs 1,476 s FedAvg	/	2,000 epochs 7,268 s FedAvg	2,000 epochs 669 s FedAvg	2,000 epochs 614 s FedAvg	/	/
default credit vertical	Tree-based model (gbdt_64_64_6)	64 trees 2,800 s SecureBoost	/	/	64 trees 9,810 s SecureBoost	/	/	/	/	64 trees 489 s SecureBoost
	Neural network (mlp_128_128_128)	1 epoch 30,952 s HE-based	/	/	/	/	/	/	10 epochs 1,354 s sMPC-based	/
<b>Communication cost</b>										
femnist horizontal	Neural network (lenet)	/	80,000 Rounds 19.71 GiB	N/A N/A	/	80,000 Rounds 19.71 GiB	80,714 Rounds 20.41 GiB	80,060 Rounds 19.96 GiB	/	/
default credit vertical	Tree-based model (gbdt_64_64_6)	2,636 Rounds N/A	/	/	50,535 Rounds 1.36 GiB	/	/	/	/	11,969 Rounds 0.39 GiB
	Neural network (mlp_128_128_128)	1,886 Rounds N/A	/	/	/	/	/	/	350,289 Rounds 69.87 GiB	/
<b>Peak memory usage in total</b>										
femnist horizontal	Neural network (lenet)	/	0.55 GiB	9.21 GiB	/	1.95 GiB	52.12 GiB	4.91 GiB	/	/
default credit vertical	Tree-based model (gbdt_64_64_6)	9.93 GiB	/	/	0.71 GiB	/	/	/	/	5.46 GiB
	Neural network (mlp_128_128_128)	17.50 GiB	/	/	/	/	/	/	0.44 GiB	/

Table 7: System performance comparison in training time, communication cost, and peak memory usage. "/" suggests the lack of functionality and "N/A" suggests missing logging due to module separation (see Appendix E). No framework consistently outperforms others in all three factors.

264 **RQ4: Which FL framework has the best system performance and what causes the differences?**  
265 Here we consider the training time, the communication cost of participants, and the peak memory  
266 consumption as metrics to evaluate system performance. In this way, the benchmark provides  
267 reference on the FL framework selection for application scenarios with different hardware and  
268 resource constraints. Table 7 shows the evaluation results.

269 We have the following observations. (1) Regarding FedAvg on femnist, Flower and FLUTE have a  
270 much lower training time than the other frameworks. FedML and TFF are slow since they launch the  
271 clients with less or no parallelism among clients' training. All frameworks have the same or close  
272 communication cost following the FedAvg algorithm. For peak memory usage, Flower has a high  
273 memory requirement as it keeps the states of all clients at all times regardless of client sampling.  
274 (2) Regarding vertical FL with the tree-based model, FedTree is significantly faster than FATE and  
275 Fedlearner. Training with the same number of trees, FATE has the lowest communication frequency  
276 and Fedlearner has the lowest peak memory usage. (3) Regarding vertical FL with neural networks,  
277 while FATE and CrypTen adopt different privacy techniques to protect the transferred messages,  
278 CrypTen is much faster than FATE with lower memory usage. However, the communication frequency  
279 of CrypTen is high. Overall, there is no framework that consistently outperforms others in all three  
280 factors (i.e., training efficiency, communication efficiency, and memory usage).

## 281 6 Discussion and future work

282 Based on our benchmark results in Section 5, here we answer the question in the introduction by  
283 providing a complete FL framework selection guideline and also discuss future works.

284 For FL practitioners to select an FL framework for a specific use case, the first step is to analyze  
285 the qualitative requirement of the use case and narrow down the scope using Table 2 and Table 3.  
286 They should also find the benchmark scenario that is most similar to their use case and refer to  
287 Table 6 for a list of preferred model types. Considering the infrastructure hardware constraint for the  
288 use case, practitioners should cross-check Table 7 and Table 2 to find frameworks that best match  
289 their deployment environment. If no existing FL framework satisfies all constraints for the use case,  
290 practitioners should consider the option of customizing one of the frameworks and can refer to Table  
291 3 to evaluate the feasibility and difficulty for further development.

292 There are the following future directions to further improve UniFed: (1) We expect more datasets  
293 can be incorporated into UniFed as the FL studies grow, especially for vertical FL. (2) We will  
294 periodically check the latest and representative FL frameworks (e.g., FedScale [29]) which is one  
295 recent open-source framework that we do not consider due to time constraints) and include them into  
296 UniFed. (3) We may discuss and evaluate the fairness and incentives of FL frameworks when there  
297 are enough frameworks enabling these factors. (4) We plan to launch an open competition to call for  
298 efficient, effective, and secure solutions using the existing FL frameworks.

## References

- 299
- 300 [1] Breast cancer wisconsin (diagnostic) data set. [https://www.kaggle.com/datasets/](https://www.kaggle.com/datasets/uciml/breast-cancer-wisconsin-data)  
301 [uciml/breast-cancer-wisconsin-data](https://www.kaggle.com/datasets/uciml/breast-cancer-wisconsin-data). Accessed: 2022-06-08.
- 302 [2] Electric motor temperature. [https://www.kaggle.com/datasets/wkirgsn/](https://www.kaggle.com/datasets/wkirgsn/electric-motor-temperature)  
303 [electric-motor-temperature](https://www.kaggle.com/datasets/wkirgsn/electric-motor-temperature). Accessed: 2022-06-08.
- 304 [3] Fedlearner. <https://github.com/bytedance/fedlearner>. Accessed: 2022-06-06.
- 305 [4] Give me some credit. <https://www.kaggle.com/c/GiveMeSomeCredit/data>. Accessed:  
306 2022-06-08.
- 307 [5] Homo secureboost. [https://github.com/FederatedAI/FATE/blob/master/doc/](https://github.com/FederatedAI/FATE/blob/master/doc/federatedml_component/ensemble.md#homo-secureboost)  
308 [federatedml\\_component/ensemble.md#homo-secureboost](https://github.com/FederatedAI/FATE/blob/master/doc/federatedml_component/ensemble.md#homo-secureboost). Accessed: 2022-06-09.
- 309 [6] Paddlefl. <https://github.com/PaddlePaddle/PaddleFL>. Accessed: 2022-06-06.
- 310 [7] Tensorflow federated. <https://github.com/tensorflow/federated>. Accessed: 2022-06-  
311 06.
- 312 [8] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving,  
313 M. Isard, et al. {TensorFlow}: A system for {Large-Scale} machine learning. In *12th USENIX*  
314 *symposium on operating systems design and implementation (OSDI 16)*, pages 265–283, 2016.
- 315 [9] D. J. Beutel, T. Topal, A. Mathur, X. Qiu, T. Parcollet, P. P. de Gusmão, and N. D. Lane. Flower:  
316 A friendly federated learning research framework. *arXiv preprint arXiv:2007.14390*, 2020.
- 317 [10] S. Caldas, S. M. K. Duddu, P. Wu, T. Li, J. Konečný, H. B. McMahan, V. Smith, and A. Talwalkar.  
318 Leaf: A benchmark for federated settings. *Workshop on Federated Learning for Data Privacy*  
319 *and Confidentiality*, 2019.
- 320 [11] A. C. Cameron, P. K. Trivedi, F. Milne, and J. Piggott. A microeconomic model of the  
321 demand for health care and health insurance in australia. *The Review of economic studies*,  
322 55(1):85–106, 1988.
- 323 [12] D. Chai, L. Wang, K. Chen, and Q. Yang. Fedeval: A benchmark system with a comprehensive  
324 evaluation model for federated learning. *arXiv preprint arXiv:2011.09655*, 2020.
- 325 [13] K. Cheng, T. Fan, Y. Jin, Y. Liu, T. Chen, D. Papadopoulos, and Q. Yang. Secureboost: A  
326 lossless federated learning framework. *IEEE Intelligent Systems*, 36(6):87–98, 2021.
- 327 [14] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik. Emnist: an extension of mnist to handwritten  
328 letters. *arXiv preprint arXiv:1702.05373*, 2017.
- 329 [15] P. Cortez and A. M. G. Silva. Using data mining to predict secondary school student performance.  
330 2008.
- 331 [16] D. Dimitriadis, M. H. Garcia, D. M. Diaz, A. Manoel, and R. Sim. Flute: A scalable, ex-  
332 tensible framework for high-performance federated learning simulations. *arXiv preprint*  
333 *arXiv:2203.13789*, 2022.
- 334 [17] D. Dua and C. Graff. UCI machine learning repository, 2017.
- 335 [18] C. Dwork, A. Roth, et al. The algorithmic foundations of differential privacy. *Found. Trends*  
336 *Theor. Comput. Sci.*, 9(3-4):211–407, 2014.
- 337 [19] A. Fallah, A. Mokhtari, and A. Ozdaglar. Personalized federated learning: A meta-learning  
338 approach. *arXiv preprint arXiv:2002.07948*, 2020.

- 339 [20] S. Hardy, W. Henecka, H. Ivey-Law, R. Nock, G. Patrini, G. Smith, and B. Thorne. Private fed-  
340 erated learning on vertically partitioned data via entity resolution and additively homomorphic  
341 encryption. *arXiv preprint arXiv:1711.10677*, 2017.
- 342 [21] C. He, K. Balasubramanian, E. Ceyani, C. Yang, H. Xie, L. Sun, L. He, L. Yang, P. S. Yu,  
343 Y. Rong, et al. Fedgraphnn: A federated learning system and benchmark for graph neural net-  
344 works. *Workshop on Distributed and Private Machine Learning: The International Conference*  
345 *on Learning Representations (DPML-ICLR)*, 2021.
- 346 [22] C. He, S. Li, J. So, X. Zeng, M. Zhang, H. Wang, X. Wang, P. Vepakomma, A. Singh, H. Qiu,  
347 et al. Fedml: A research library and benchmark for federated machine learning. *arXiv preprint*  
348 *arXiv:2007.13518*, 2020.
- 349 [23] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In  
350 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–  
351 778, 2016.
- 352 [24] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–  
353 1780, 1997.
- 354 [25] S. Hu, Y. Li, X. Liu, Q. Li, Z. Wu, and B. He. The oarf benchmark suite: Characterization and  
355 implications for federated learning systems. *arXiv preprint arXiv:2006.07856*, 2020.
- 356 [26] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz,  
357 Z. Charles, G. Cormode, R. Cummings, et al. Advances and open problems in federated learning.  
358 *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- 359 [27] B. Knott, S. Venkataraman, A. Hannun, S. Sengupta, M. Ibrahim, and L. van der Maaten.  
360 CrypTen: Secure multi-party computation meets machine learning. In *arXiv 2109.00984*, 2021.
- 361 [28] B. Knott, S. Venkataraman, A. Hannun, S. Sengupta, M. Ibrahim, and L. van der Maaten.  
362 CrypTen: Secure multi-party computation meets machine learning. *Advances in Neural Informa-*  
363 *tion Processing Systems*, 34, 2021.
- 364 [29] F. Lai, Y. Dai, S. S. Singapuram, J. Liu, X. Zhu, H. V. Madhyastha, and M. Chowdhury.  
365 FedScale: Benchmarking model and system performance of federated learning. In *International*  
366 *Conference on Machine Learning (ICML)*, 2022.
- 367 [30] Y. LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist>.  
368 Accessed: 2022-06-08.
- 369 [31] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D.  
370 Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*,  
371 1(4):541–551, 1989.
- 372 [32] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document  
373 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 374 [33] O. Li, J. Sun, X. Yang, W. Gao, H. Zhang, J. Xie, V. Smith, and C. Wang. Label leakage and  
375 protection in two-party split learning. In *International Conference on Learning Representations*,  
376 2022.
- 377 [34] Q. Li, Y. Cai, Y. Han, C. M. Yung, T. Fu, and B. He. Fedtree: A fast, effective, and secure  
378 tree-based federated learning system. [https://github.com/Xtra-Computing/FedTree/  
379 blob/main/FedTree\\_draft\\_paper.pdf](https://github.com/Xtra-Computing/FedTree/blob/main/FedTree_draft_paper.pdf), 2022.
- 380 [35] Q. Li, Y. Diao, Q. Chen, and B. He. Federated learning on non-iid data silos: An experimental  
381 study. In *IEEE International Conference on Data Engineering*, 2022.

- 382 [36] Y. Liang, Y. Guo, Y. Gong, C. Luo, J. Zhan, and Y. Huang. Flbench: A benchmark suite for  
383 federated learning. In *BenchCouncil International Federated Intelligent Computing and Block  
384 Chain Conferences*, pages 166–176. Springer, 2020.
- 385 [37] B. Y. Lin, C. He, Z. Zeng, H. Wang, Y. Huang, M. Soltanolkotabi, X. Ren, and S. Avestimehr.  
386 Fednlp: A research platform for federated learning in natural language processing. *NAACL  
387 Findings*, 2022.
- 388 [38] Y. Liu, T. Fan, T. Chen, Q. Xu, and Q. Yang. Fate: An industrial grade platform for collaborative  
389 learning with data protection. *Journal of Machine Learning Research*, 22(226):1–6, 2021.
- 390 [39] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings  
391 of the IEEE international conference on computer vision*, pages 3730–3738, 2015.
- 392 [40] M. Luo, F. Chen, D. Hu, Y. Zhang, J. Liang, and J. Feng. No fear of heterogeneity: Classifier  
393 calibration for federated learning with non-iid data. *Advances in Neural Information Processing  
394 Systems*, 34, 2021.
- 395 [41] Y. Ma, D. Yu, T. Wu, and H. Wang. Paddlepaddle: An open-source deep learning platform from  
396 industrial practice. *Frontiers of Data and Computing*, 1(1):105–115, 2019.
- 397 [42] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient  
398 learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages  
399 1273–1282. PMLR, 2017.
- 400 [43] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,  
401 N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning  
402 library. *Advances in neural information processing systems*, 32, 2019.
- 403 [44] J. Siebert. *Vehicle Recognition Using Rule Based Methods*. Turing Institute, 1987.
- 404 [45] Y. Wu, S. Cai, X. Xiao, G. Chen, and B. C. Ooi. Privacy preserving vertical federated learning  
405 for tree-based models. *arXiv preprint arXiv:2008.06170*, 2020.
- 406 [46] A. C.-C. Yao. How to generate and exchange secrets. In *27th Annual Symposium on Foundations  
407 of Computer Science (sfcs 1986)*, pages 162–167. IEEE, 1986.
- 408 [47] I.-C. Yeh and C.-h. Lien. The comparisons of data mining techniques for the predictive accuracy  
409 of probability of default of credit card clients. *Expert systems with applications*, 36(2):2473–  
410 2480, 2009.
- 411 [48] Q. Zhang, C. Wang, H. Wu, C. Xin, and T. V. Phuong. Gelu-net: A globally encrypted, locally  
412 unencrypted deep neural network for privacy-preserved learning. In *IJCAI*, pages 3933–3939,  
413 2018.
- 414 [49] Y. Zhang and H. Zhu. Additively homomorphical encryption based deep neural network for  
415 asymmetrically collaborative machine learning. *arXiv preprint arXiv:2007.06849*, 2020.
- 416 [50] B. Zhu, L. Wang, Q. Pang, S. Wang, J. Jiao, D. Song, and M. I. Jordan. Byzantine-robust  
417 federated learning with optimal statistical rates and privacy guarantees. *arXiv preprint  
418 arXiv:2205.11765*, 2022.

419 **Checklist**

- 420 1. For all authors...
- 421 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
422 contributions and scope? [Yes]
- 423 (b) Did you describe the limitations of your work? [Yes] We discuss our limitations and  
424 the future work in Section 6
- 425 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We  
426 provide a benchmark for existing FL frameworks so it is not closely relevant.
- 427 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
428 them? [Yes]
- 429 2. If you are including theoretical results...
- 430 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 431 (b) Did you include complete proofs of all theoretical results? [N/A]
- 432 3. If you ran experiments (e.g. for benchmarks)...
- 433 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
434 imental results (either in the supplemental material or as a URL)? [Yes] We include  
435 those in the supplemental material and provide a GitHub link for our implementation.
- 436 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
437 were chosen)? [Yes] We include them in the supplemental material and GitHub  
438 implementation.
- 439 (c) Did you report error bars (e.g., with respect to the random seed after running ex-  
440 periments multiple times)? [No] We observe different stability properties for FL  
441 frameworks with different ML backends, which makes it inconsistent to report error  
442 bars across all frameworks. However, we verified that, the conclusions we present in  
443 the paper are consistent in different runs with different random seeds, which can be  
444 reproduced with the codebase we provide. (See Section 5)
- 445 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
446 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 5.1
- 447 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 448 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 449 (b) Did you mention the license of the assets? [Yes] It is mentioned in the cited source.
- 450 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 451 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
452 using/curating? [N/A] We only use existing datasets in our work.
- 453 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
454 information or offensive content? [N/A] We only use existing datasets in our work.
- 455 5. If you used crowdsourcing or conducted research with human subjects...
- 456 (a) Did you include the full text of instructions given to participants and screenshots, if  
457 applicable? [N/A]
- 458 (b) Did you describe any potential participant risks, with links to Institutional Review  
459 Board (IRB) approvals, if applicable? [N/A]
- 460 (c) Did you include the estimated hourly wage paid to participants and the total amount  
461 spent on participant compensation? [N/A]