

FedScale: Benchmarking Model and System Performance of Federated Learning

Fan Lai, Yinwei Dai, Xiangfeng Zhu, Mosharaf Chowdhury

University of Michigan

Abstract

1 We present FedScale, a diverse set of challenging and realistic benchmark datasets
2 to facilitate scalable, comprehensive, and reproducible federated learning (FL)
3 research. FedScale datasets are large-scale, encompassing a diverse range of im-
4 portant FL tasks, such as image classification, object detection, language modeling,
5 speech recognition, and reinforcement learning. For each dataset, we provide a
6 unified evaluation protocol using realistic data splits and evaluation metrics. To
7 meet the pressing need for reproducing realistic FL at scale, we have also built
8 an efficient evaluation platform to simplify and standardize the process of FL ex-
9 perimental setup and model evaluation. Our evaluation platform provides flexible
10 APIs to implement new FL algorithms and includes new execution backends with
11 minimal developer efforts. Finally, we perform indepth benchmark experiments
12 on these datasets. Our experiments suggest fruitful opportunities in heterogeneity-
13 aware co-optimizations of the system and statistical efficiency under realistic FL
14 characteristics. FedScale is open-source with permissive licenses and actively
15 maintained,¹ and we welcome feedback and contributions from the community.

16 1 Introduction

17 Federated learning (FL) is an emerging machine learning (ML) setting where a logically
18 centralized coordinator orchestrates many distributed clients (e.g., smartphones or laptops)
19 to collaboratively train or evaluate a model [14, 32] (Figure 1). It enables model training
20 and evaluation on end-user data, while circumventing high cost and privacy risks in
21 gathering the raw data from clients, with applications in diverse domains: for example,
22 NVIDIA applies FL to create medical imaging AI [38]; Google runs federated training
23 of NLP models in Google keyboard [17, 55];
24 Apple performs federated evaluation and tuning of automatic speech recognition models on end-user
25 devices [43]; IBM is deploying FL infrastructure to help detect financial misconducts [39].

32 To address challenges arising from the heterogeneous execution speeds of client devices as well
33 as non-IID data distributions, existing efforts have focused on optimizing different aspects of FL:
34 (1) *System efficiency*: reducing computation load (e.g., using smaller models [47]) or communication
35 traffic (e.g., local SGD [42]) to achieve faster on-device execution; (2) *Statistical efficiency*: designing

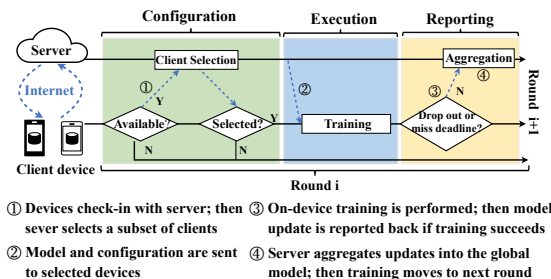


Figure 1: Standard FL protocol [14, 54].

¹FedScale is available at <https://github.com/SymbioticLab/FedScale>.

Features	OARF [30]	LEAF [15]	FedEval [16]	FedML [28]	Flower [13]	FedScale
Heter. Client Dataset	○	○	✗	○	○	✓
Heter. System Speed	✗	✗	✗	✗	✗	✓
Client Availability	✗	✗	✗	✗	✗	✓
Scalable Platform	✗	✗	○	○	✓	✓
Real FL Runtime	✗	✗	✗	✗	✗	✓
Flexible APIs	✗	✗	✗	✓	✓	✓

Table 1: Comparing FedScale with existing FL benchmarks and libraries. ○ implies limited support. We provide more details for this comparison in Appendix B.

36 data heterogeneity-aware algorithms (e.g., client clustering [26]) to obtain better training accuracy
 37 with fewer training rounds; (3) *Privacy and security*: developing reliable strategies (e.g., differentially
 38 private training [31]) to make FL more privacy-preserving and robust to potential attacks.

39 The performance of an FL solution greatly depends on the characteristics of data, device capabilities,
 40 and participation of clients; overlooking any one aspect can mislead FL evaluation (§2). For example,
 41 dynamics of client system performance or availability (e.g., device drop-out or rejoining) can affect
 42 the dynamics of data availability (distribution shift of cross-device data), which may impair model
 43 convergence [20]; too few clients can lead to unstable statistical training convergence, but too many
 44 can slow down practical model aggregation because of heterogeneous system speed. As such, a
 45 comprehensive suite of benchmarks that combine diverse aspects of practical FL is crucial for
 46 systemic evaluation and comparison of different efforts.

47 Existing benchmarks for FL are mostly borrowed from traditional ML benchmarks (e.g., MLPerf [40])
 48 or designed for simulated FL environments like TensorFlow Federated [12] or PySyft [8]. As shown
 49 in Table 1, existing benchmarks for FL fall short in multiple ways: (1) they are limited in the versatility
 50 of data for various real-world FL applications. Indeed, even though they may have quite a few datasets
 51 and FL training tasks (e.g., FedEval [16] and LEAF [15]), their datasets often contain synthetically
 52 generated partitions derived from conventional datasets (e.g., CIFAR) and do not represent realistic
 53 characteristics; (2) existing benchmarks often overlook different aspects of practical FL. For example,
 54 system speed and availability of the client are largely missing (e.g., FedML [8] and Flower [13]),
 55 which discourages FL efforts from considering system efficiency and leads to overly optimistic
 56 statistical performance (§2); (3) their experimental environments are unable to reproduce the practical
 57 scale of FL deployments. While real FL often involves thousands of participants in each training
 58 round [32, 55], existing benchmarking platforms – therefore, many existing FL solutions – are merely
 59 able to support the training of tens of participants per round; (4) they may lack user-friendly APIs
 60 for automated integration, resulting in great engineering efforts in benchmarking new plugins.

61 **Contributions:** In this paper, we introduce FedScale, an FL benchmark to empower comprehensive
 62 and standardized FL evaluations. As shown in Figure 2,
 63 we make the following contributions:

- 64 • To the best of our knowledge, we incorporate the most comprehensive FL datasets for evaluating different
 65 aspects of real FL deployments. FedScale currently
 66 has 18 realistic FL datasets spanning across small,
 67 medium, and large scales for a wide variety of task cat-
 68 egories, such as image classification, object detection,
 69 language modeling, speech recognition, machine trans-
 70 lation, recommendation, and reinforcement learning.
 71 To account for practical client behaviors, we include
 72 real-world measurements of mobile devices, and asso-
 73 ciate each client with his computation and communi-
 74 cation speeds, as well as availability status dynamics.
- 75 • We build an automated evaluation platform, FedScale
 76 Automated Runtime (FAR), to simplify and standard-
 77 ize the FL evaluation in a more realistic setting. FAR
 78 integrates real FL statistical and system dataset, and
 79

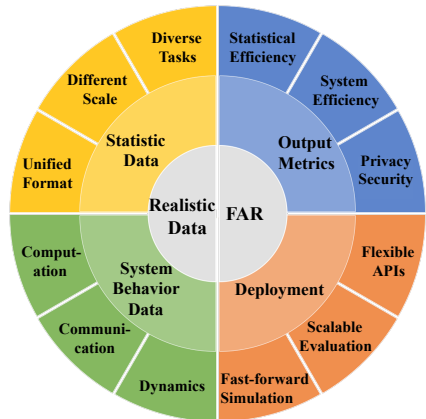


Figure 2: FedScale provides real FL data and an automated evaluation platform.

80 thus can pinpoint various practical FL metrics needed in today’s work. FAR allows easy deploy-
 81 ment of new plugins with flexible APIs and can perform the training of thousands of clients in
 82 each round on a few GPUs efficiently. FAR is built atop of our recent work Oort [36], which has
 83 passed a rigorous artifact evaluation in OSDI 2021.

- 84 • We perform indepth benchmark experiments for recent FL efforts in FedScale setting, and
 85 highlight the pressing need of co-optimizing system and statistical efficiency in a heterogeneity-
 86 aware manner, especially in tackling system stragglers and biased model performance.

87 2 Background

88 **Existing efforts optimize for various goals of practical FL** To tackle heterogeneous client data,
 89 FedProx [37], FedYogi [44] and Scaffold [33] introduce adaptive client/server optimizations that use
 90 control variates to correct for the ‘drift’ in model updates. Instead of training a single global model,
 91 some efforts resort to training a mixture of models [19, 22], clustering clients over training [27],
 92 or enforcing guided client selection [36]; To tackle the scarce and heterogeneous device resource,
 93 FedAvg [42] reduces communication cost by performing multiple local SGD steps, while some
 94 works compress the model update by filtering out or quantizing unimportant parameters [46, 34];
 95 After realizing the privacy risk in FL [24, 51], DP-SGD [25] enhances the privacy by introducing
 96 differential privacy, and DP-FTRL [31] applies the tree aggregation to add noise to the sum of
 97 mini-batch gradients to ensure privacy further. These FL efforts often navigate privacy-accuracy-
 98 computation trade-offs. As such, a realistic FL setting is crucial for comprehensive evaluations.

99 **Existing FL benchmarks can be misleading** Existing benchmarks often lack realistic client
 100 statistical and system behavior datasets, and/or fail to reproduce large-scale FL deployments.
 101 Unfortunately, these limitations imply that
 102 they are not only insufficient for bench-
 103 marking diverse FL optimizations, but they
 104 can even mislead performance evaluations:
 105 (1) As shown in Figure 3(a), the statistical
 106 performance becomes worse when encoun-
 107 tering practical client behaviors (e.g., strag-
 108 glers and training failures), which indicates
 109 that existing benchmarks that do not have
 110 systems traces can produce overly opti-
 111 mistic statistical performance by overlook-
 112 ing systems characteristics; (2) FL training
 113 with hundreds of participants each round
 114 performs better than that with tens of par-
 115 ticipants (Figure 3(b)). As such, existing benchmark platforms can under-report existing FL optimiza-
 116 tions as they cannot support the practical FL scale with a large number of participants.

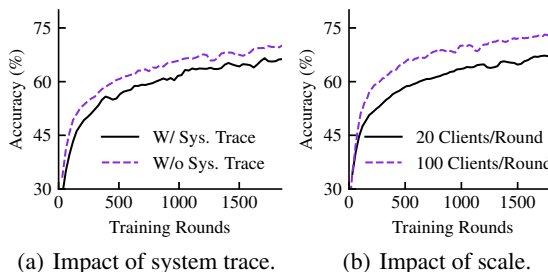


Figure 3: Existing benchmarks can mislead.²

117 3 FedScale Dataset: Realistic Workloads for Federated Learning

118 FL performance relies on at least three aspects: (1) *Client statistical data*: the client dataset used
 119 for training or testing determines the statistical efficiency of FL tasks (e.g., convergence and model
 120 accuracy); (2) *Client system behavior*: the compute/communication speed of the client device and its
 121 availability over time determine the system efficiency of FL tasks (e.g., duration of each round and
 122 physical cost) and the availability of statistical data; and (3) *Task categories*: model and application
 123 combinations that are running can exhibit different reliance on client statistical data and execute at
 124 different system speeds. Because client data is tightly coupled with the client device, these aspects
 125 interplay with each other and can impact the performance of an FL optimization, be it for statistical
 126 efficiency, system efficiency, or privacy. As such, an ideal suite of FL benchmarking dataset should
 127 cover all three aspects and support FL deployments at diverse scales.

128 We next introduce how we collected and partitioned realistic datasets in order to generate a versatile
 129 suite of FL datasets provided in FedScale.

²We train the ShuffleNet model on OpenImage classification task. More experimental setups in Section 5.

Category	Name	Data Type	#Clients	#Instances	Example Task
CV	iNature	Image	2,295	193K	Classification
	OpenImage	Image	13,771	1.3M	Classification, Object detection
	Google Landmark	Image	43,484	3.6M	Classification
	Charades	Video	266	10K	Action recognition
	VLOG	Video	4,900	9.6K	Classification, Object detection
	Waymo Motion	Video	496,358	32.5M	Motion prediction
NLP	Europarl	Text	27,835	1.2M	Text translation
	Blog Corpus	Text	19,320	137M	Word prediction
	Reddit	Text	1,660,820	351M	Word prediction
	CoQA	Text	7,189	114K	Question Answering
	LibriTTS	Text	2,456	37K	Text to speech
	Google Speech	Audio	2,618	105K	Speech recognition
	Common Voice	Audio	12,976	1.1M	Speech recognition
Misc ML	Taobao	Text	182,806	20.9M	Recommendation
	Fox Go	Text	150,333	4.9M	Reinforcement learning

Table 2: Statistics of *partial* FedScale datasets (the full list and more details of data and its partition are in Appendix A). FedScale has 18 real-world federated datasets; each dataset is partitioned by its real client-data mapping. Note that we remove the sensitive information in these datasets.

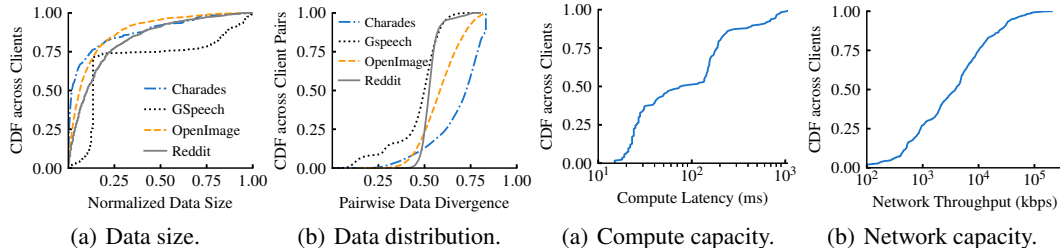


Figure 4: Non-IID client data distribution.

Figure 5: Heterogeneous client system speed.

130 3.1 Client Statistical Dataset

131 FedScale currently has 18 realistic FL datasets (Table 2), which can be used in various FL tasks (e.g.,
 132 federated training/testing or on-device fine-tuning). The raw data of these datasets are collected from
 133 different sources and stored in various formats. We clean up the raw data, partition them into new FL
 134 datasets, and streamline new datasets into consistent formats. Moreover, we categorize them into
 135 different FL use cases and provide Python APIs for integrating them into today’s frameworks.

136 **Realistic data and partitions** We target realistic datasets with client information, and partition the
 137 raw dataset using the unique client identification. For example, OpenImage is a vision dataset collected
 138 by Flickr, wherein different mobile users upload their images to the cloud for public use. We use
 139 the AuthorProfileUrl attribute of the OpenImage data to map data instances to each client, whereby
 140 we extract the realistic distribution of the raw data. Following existing FL deployments [55], for
 141 each dataset, we assign its clients into the training, validation or testing groups, whereby we get the
 142 training, validation and testing set for it. Here, we pick four real-world datasets – video (Charades),
 143 audio (Google Speech), image (OpenImage), and text (Reddit) – to illustrate the characteristics of FL.
 144 Each dataset consists of hundreds or up to millions of clients and millions of data points. Figure 4
 145 reports the *Cumulative Distribution Function* (CDF) of the data distribution, wherein we see a high
 146 statistical deviation across clients not only in the quantity of samples (Figure 4(a)) but also in the data
 147 distribution (Figure 4(b)).³ Our findings confirm the non-IID data distribution in FL.

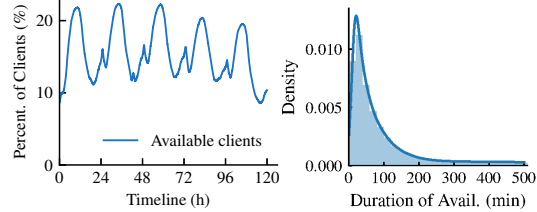
³We report the pairwise Jensen–Shannon distance of the categorical distribution between two clients.

148 **Different scales across diverse task categories** To accommodate diverse scenarios in practical
 149 FL, FedScale includes small-, medium-, and large-scale datasets across a wide range of tasks, from
 150 hundreds to millions of clients. Some datasets can be applied in different tasks, as we enrich their use
 151 case by driving different metadata from the same raw data. For example, the raw OpenImage dataset
 152 can be used for object detection, and we extract each object therein and generate a new dataset for
 153 image classification. Moreover, we provide APIs for the developer to customize their dataset (e.g.,
 154 enforcing new data distribution or extracting a subset of clients).

155 3.2 Client System Behavior Dataset

156 **Client device system speed is heterogeneous** We formulate the system trace of different clients
 157 using *AI Benchmark* [1] and *MobiPerf Measurements* [7] on mobiles. *AI Benchmark* provides the
 158 training and inference speed of diverse models (e.g., MobileNet) across a wide range of device
 159 models (e.g., Huawei P40 and Samsung Galaxy S20), while *MobiPerf* has collected the available
 160 cloud-to-edge network throughput of over 100k world-wide mobile clients. As specified in real
 161 FL deployments [14, 55], we focus on mobile devices that have larger than 2GB RAM and connect
 162 with WiFi; Figure 5 reports that their compute and network capacity can exhibit order-of-magnitude
 163 difference. As such, how to orchestrate scarce resources and mitigate stragglers are paramount for
 164 high system efficiency.

165 **Client device availability is dynamic** We in-
 166 corporate a large-scale user behavior dataset
 167 spanning 136k users [54] to emulate the behav-
 168 iors of clients. It includes 180 million trace
 169 items of client devices (e.g., battery charge or
 170 screen lock) over a week. We follow the real FL
 171 setting, which considers the device in charging
 172 to be available [12] and observe great dynamics
 173 in their availability: (i) the number of available
 174 clients reports diurnal variation (Figure 6(a)).
 175 This confirms the cyclic patterns in the client
 176 data, which can deteriorate the statistical per-
 177 formance of FL [20]. (ii) the duration of each
 178 available slot is not long-lasting (Figure 6(b)).



(a) Inter-device availability. (b) Intra-device availability.

Figure 6: Client availability is dynamic.

179 This highlights the need of handling failures (clients become offline) during training, since the
 180 duration of each round (also a number of minutes) is comparable to that of each available slot.

181 4 FAR: Evaluation Platform for Federated Learning

182 Existing FL evaluation platforms can
 183 hardly reproduce the scale of practical FL
 184 deployments and fall short in providing
 185 user-friendly APIs, which requires great de-
 186 veloper efforts to deploy new plugins. As
 187 such, we introduce FedScale Automated
 188 Runtime (FAR), an automated and easily-
 189 deployable evaluation platform, to simplify
 190 and standardize the FL evaluation under
 191 a practical setting. FAR is based on our
 192 Oort project [36], which has been shown to
 193 scale well and can emulate FL training of
 194 thousands of clients in each round.

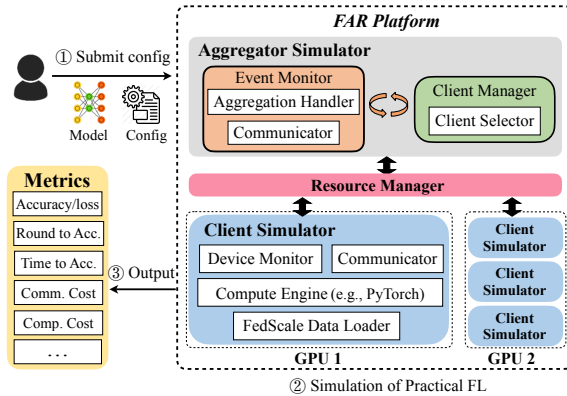


Figure 7: FAR enables the developer to benchmark various FL efforts with practical FL data and metrics.

195 **Overview of FedScale Automated Run-
 196 time (FAR)** FAR is an automated evalua-
 197 tion platform that can emulate realistic FL
 198 behaviors on GPU/CPU, while providing

Module	API Name	Example Use Case
Aggregator Simulator	<code>round_initialization_handler(*args)</code>	Client clustering
	<code>round_completion_handler(*args)</code>	Adaptive/secure model aggregation
	<code>client_completion_handler(client_id, msg)</code>	Straggler mitigation
	<code>push_msg_to_client(client_id, msg)</code>	Model compression
Client Manager	<code>select_clients(*args)</code>	Client selection
	<code>select_model_for_client(client_id)</code>	Adaptive model selection
Client Simulator	<code>train(client_data, model, config)</code>	Local SGD/malicious attack
	<code>push_msg_to_aggregator(msg)</code>	Model compression

Table 3: Some example APIs. FedScale provides APIs to deploy new plugins for various designs.

199 various practical FL metrics, such as computation/communication cost, latency and wall clock time,
200 for evaluating today’s efforts. As shown in Figure 7, FAR primarily consists of three components:

- 201 • *Aggregator Simulator*: It acts as the aggregator in practical FL, which selects participants,
202 distributes execution profiles (e.g., model weight), and handles result (e.g., model updates)
203 aggregation. In each round, its client manager uses the client behavior trace to decide whether a
204 client is available; then it selects the specified number of clients to participate that round. Once
205 receiving new events, the event monitor will activate the handler (e.g., aggregation handler to
206 perform model aggregation), or the communicator to send/receive messages. The communicator
207 records the size (cost) of every network traffic, and its FL runtime latency ($\frac{traffic_size}{client_bandwidth}$).
- 208 • *Client Simulator*: It works as the client in FL. FedScale data loader loads the federated dataset of
209 that client and feeds this data to the compute engine to run real training/testing. The computation
210 latency is determined by ($\#_processed_sample \times latency_per_sample$), and the communica-
211 tor handles the network traffics and records the communication latency ($\frac{traffic_size}{client_bandwidth}$). At the
212 same time, the device monitor handles different function calls specified by the developer; it will
213 also terminate the simulation of this client and report failure(s) if the current runtime exceeds the
214 available slot (indicated in the client availability trace).
- 215 • *Resource Manager*: It orchestrates the available physical resource for evaluation to maximize the
216 utilization of resource. For example, when the number of participants in that round exceeds the
217 resource capacity (e.g., simulating thousands of clients on a few GPUs), the resource manager
218 queues the overcommitted tasks of clients and schedules a new client simulation request from
219 this queue once resource becomes available.

220 Note that capturing runtime performance (e.g., wall clock time of training) is rather slow in practical
221 FL (each client takes several minutes), but FAR enables *fast-forward* simulation for interactive
222 development, since the real training on our platform often takes only a few seconds per round.

223 **FAR enables automated and standardized FL simulation** FAR incorporates realistic FL traces,
224 using the aforementioned trace by default, to automatically emulate the practical FL workflow: ①
225 *Task submission*: FL developers specify their configurations (e.g., model and dataset), which can
226 be federated training or testing, and the FAR resource manager will initiate the aggregator and
227 client simulator on available resource (GPU, CPU, other accelerators, or even smartphones); ② *FL*
228 *simulation*: This evaluation stage follows the standardized FL lifecycle (in Figure 1). In each training
229 round, the aggregator inquires the client manager to select the participants, whereby the resource
230 manager distributes the client configuration to the available client simulators. After the completion
231 of each client, the client simulator pushes the model update to the aggregator, which then performs
232 the model aggregation. ③ *Metrics output*: During training, the developer can query the practical
233 evaluation metrics on the fly. Figure 7 lists some popular metrics supported in FAR.

234 **FAR is easily-deployable and extensible for plugins** FAR provides flexible APIs, which can
235 accommodate with different execution backends (e.g., PyTorch and TensorFlow) by design, for the
236 developer to quickly deploy new plugins for customized evaluations. Table 3 illustrates some example
237 APIs that can facilitate diverse FL efforts, and Figure 9 dictates an example showing how these APIs

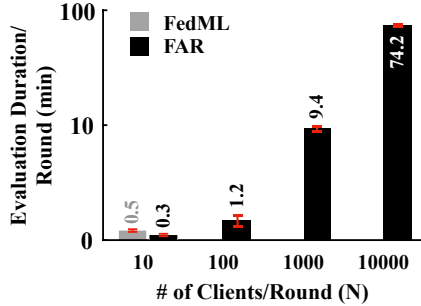


Figure 8: FAR can support thousands of clients per round, while FedML failed to run even 100 clients.

```

from fedscale.core.client import Client

class Customized_Client(Client):
    # Redefine training (e.g., for local
    # SGD/gradient compression)
    def train(self, client_data, model, conf):
        # Code of plugin
        ...
        # Results will be sent to aggregator
        return training_result

```

Figure 9: Add plugins by inheritance.

238 help to benchmark a new design of local client training with a few lines of code. Specifically, the
 239 developer can redefine client training function `run_client` by inheriting the base `Executor` module,
 240 and this plugin will be automatically integrated into FedScale during evaluations. Moreover, FAR
 241 can embrace new realistic (statistical client or system behavior) datasets with the built-in APIs. For
 242 example, the developer can import his own dataset of the client availability by leveraging the API
 243 (`load_client_availability`), and FAR will automatically force this trace during evaluations. We
 244 also provide more examples in Appendix C to demonstrate the ease of evaluating different today’s FL
 245 algorithms in FAR— a few lines are all we need!

246 **FAR is scalable and efficient** FAR can perform large-scale simulations (e.g., thousands of partici-
 247 pants in each round) in both standalone (single CPU/GPU) and distributed (multiple CPUs/GPUs)
 248 setting. This is because: (1) FAR can support multiprocessing on a single GPU so that multiple
 249 client simulators can co-locate on the same GPU; (2) [our resource manager monitors the fine-grained](#)
 250 [resource utilization of machines, queues the overcommitted simulation requests, adaptively dispatches](#)
 251 [simulation requests of the client across machines to achieve load balance, and then orchestrates the](#)
 252 [simulation based on the client mirror clock](#); (3) FAR maximizes the resource utilization by overlap-
 253 [ping the communication and computation phrases of different clients](#). For example, the simulator can
 254 [turn to train new clients while the communication of the last client is on the fly](#). As shown in Figure 8
 255 ⁴, FAR not only runs faster than FedML [28] (using 10 clients per round), thus saving lots of GPU
 256 hours, but can support large-scale evaluations efficiently. Instead, state-of-the-art platforms hardly
 257 support the practical FL scale with hundreds of clients, because they mostly rely on the traditional
 258 ML architectures (e.g., the primitive parameter-server architecture), which are primarily designed for
 259 the traditional ML training on a number of workers with large batch size.

260 5 Experiments

261 In this section, we first show how FedScale can benefit the benchmarking of existing efforts optimizing
 262 for different aspects of FL. Moreover, we highlight some important insights to improve practical FL.

263 **Experimental setup** We use 10 NVIDIA Tesla P100 GPUs in our evaluations. Following the real
 264 FL deployments [14, 55], the aggregator collects updates from the first N completed participants
 265 out of $1.3N$ participants to mitigate system stragglers in each round, and $N = 100$ by default. We
 266 pick two representative datasets in FedScale, which belong to different scales and tasks: (1) *Speech*
 267 *Recognition*: the small-scale Google Speech dataset, with 105K speech commands over 2600 clients.
 268 We train ResNet-18 [29] to recognize the command among 35 categories. (2) *Image Classification*:
 269 the middle-scale OpenImage dataset, with 1.3M images spanning 600 categories across 14k clients.
 270 We train ShuffleNet-V2 [57] to classify the image. These applications and models are widely used on
 271 mobile devices. We set the minibatch size of each participant to 20, and the number of local steps K
 272 to 20. We cherry-pick the hyper-parameters with grid search, ending up with an initial learning rate
 273 0.04. These settings are consistent with the literature.

⁴We train the ShuffleNet model on OpenImage classification task. More experimental setups in Section 5.

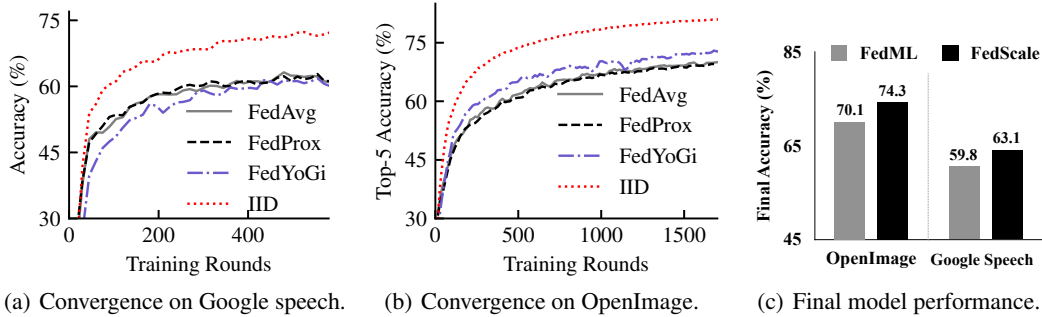


Figure 10: FedScale can benchmark the statistical FL performance. (c) shows existing benchmarks can under-report the FedYoGi performance as they cannot support a large number of participants.

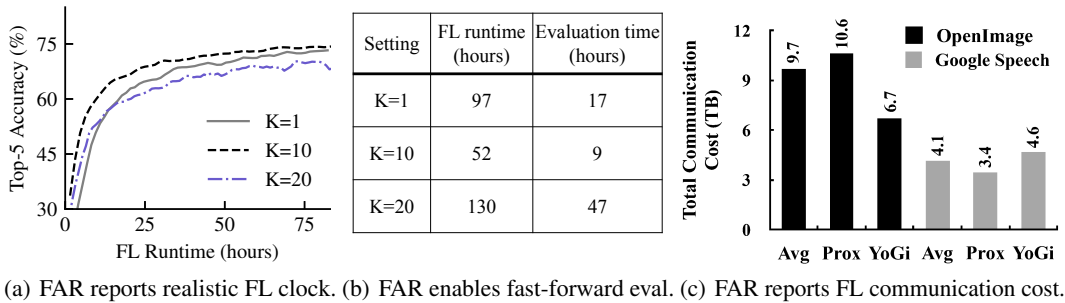


Figure 11: FedScale can benchmark realistic FL runtime. (a) and (b) report FedYoGi results on OpenImage with different number of local steps (K); (b) reports the FL runtime to reach convergence.

274 5.1 How Does FedScale Help FL Benchmarking?

275 Existing benchmarks are insufficient to evaluate the various metrics needed in today’s FL, and can even
 276 mis-report the FL performance due to their inability to reproduce the FL setting. Next, we crystallize
 277 the effectiveness of FedScale in benchmarking the different FL aspects over its counterparts.

278 **Benchmarking FL statistical efficiency.** FedScale provides various realistic client datasets to
 279 benchmark the statistical efficiency of FL optimizations. Here, we experiment with three state-of-
 280 the-art optimizations (FedAvg, FedProx and FedYoGi) – each reinvents local SGD to mitigate the
 281 data heterogeneity – and the traditional IID data setting. Figure 10 reports the statistical training
 282 convergence, and we observe that: (1) while the round-to-accuracy performance and final model
 283 accuracy of non-IID settings are consistently worse than that of the IID setting, different tasks
 284 can have different preferences on the optimizations. For example, FedYoGi performs the best on
 285 OpenImage, but it is inferior to FedAvg on Google Speech. Existing benchmarks, however, are
 286 limited to quite a few FL tasks and scales, which can discourage the evaluation of FL efforts; and (2)
 287 existing benchmarks can under-report the FL performance due to their inability to reproduce the FL
 288 setting. Figure 10(c) reports the final model accuracy using FedML and FedScale, where we attempt
 289 to reproduce the scale of practical FL with 100 participants per round in both frameworks, but FedML
 290 can only support 30 participants because of its suboptimal scalability. We notice this inability of
 291 existing benchmarks caps the practical FL performance that the algorithm can indeed achieve.

292 **Benchmarking FL system efficiency.** Existing system optimizations for FL focus on the practical
 293 runtime (e.g., wall-clock time in real FL training) and the FL execution cost. Unfortunately, existing
 294 benchmarks can hardly evaluate the FL runtime due to the lack of realistic system traces, but we now
 295 show how FedScale can help such benchmarking: (1) FAR enables fast-forward evaluations of the
 296 practical FL wall-clock time with fewer evaluation hours. Taking different number of local steps K
 297 in local SGD as an example [42], Figure 11(a) and Table 11(b) illustrate that FedScale can evaluate this

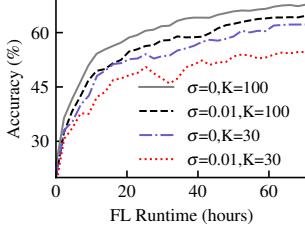


Figure 12: FedScale can benchmark privacy efforts in more realistic FL settings.

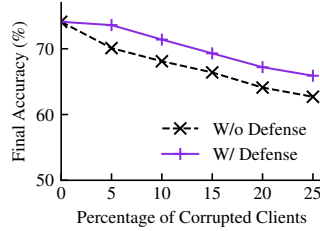


Figure 13: FedScale can benchmark security optimizations with realistic FL data.

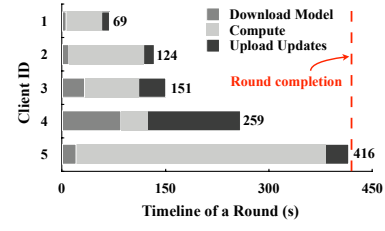


Figure 14: System stragglers greatly slow down model aggregation in practical FL.

298 impact of K on practical FL runtime in a few hours. This allows the developer to evaluate large-scale
 299 system optimizations efficiently; and (2) FAR can dictate the FL execution cost by using realistic
 300 system traces. For example, Figure 11(c) reports the practical FL communication cost in achieving
 301 the performance of Figure 10, while Figure 14 reports the system duration of individual clients. These
 302 system metrics can facilitate developers to navigate the accuracy-cost trade-off.

303 **Benchmarking FL privacy and security.** FedScale can evaluate the statistical and system effi-
 304 ciency for privacy and security optimizations in more realistic FL settings than its counterparts. Here,
 305 we give an example of how FedScale can benchmark the DP-SGD [25, 31], which applies differential
 306 privacy to improve the client privacy. We experiment with different privacy target σ ($\sigma=0$ indicates
 307 no privacy enhancement) and different number of participants per round K . Figure 12 shows that the
 308 current scale of participants (e.g., $K=30$) that today’s benchmarks can support can mislead the privacy
 309 evaluations too: while we notice great performance degradation in the training convergence of taking
 310 the privacy optimization (i.e., $\sigma=0.01$) when $K=30$, this performance drop is decent in the practical
 311 FL scale ($K=100$). Instead, FedScale is able to benchmark their performance in more FL realistic
 312 settings for various privacy use cases, such as wall-clock time, communication cost introduced in the
 313 privacy optimization, and the number of rounds needed to leak the privacy on realistic client data.

314 As for benchmarking the FL security, we follow the example setting of recent backdoor attacks [50,
 315 51] on the OpenImage, where corrupted clients flip their ground-truth labels to poison the training.
 316 We benchmarked two settings: one without security enhancement, while the other one clips the
 317 model updates as [50]. As shown in Figure 13, state-of-the-art optimizations can mitigate the attacks
 318 without hurting the overall performance when a small fraction of clients are corrupted. However,
 319 more enhancements are needed as we notice a great accuracy drop as more clients become corrupted.

320 5.2 Opportunities for Future FL Optimizations

321 **Heterogeneity-aware co-optimizations of communication and computation** Existing opti-
 322 mizations for the system efficiency often apply the same strategy on all clients (e.g., us-
 323 ing the same number of local steps [42] or compression threshold [46]), while ignoring the
 324 heterogeneous client system speed. When we outline the timeline of 5 randomly picked
 325 participants in our training of the ShuffleNet (Figure 14), we
 326 find that: (1) system stragglers can greatly slow down the round
 327 aggregation in practical FL; and (2) simply optimizing the com-
 328 munication or computation efficiency may not lead to faster
 329 rounds, as the last participant can be bottlenecked by the other
 330 resource. Here, optimizing the communication can greatly ben-
 331 efit *Client 4*, but it achieves marginal improvement on the round
 332 duration as *Client 5* is bottlenecked by computation. As such,
 333 there is an urgent need of co-optimizing the communication
 334 and computation efficiency while being heterogeneity-aware.

335 **Co-optimizations of statistical and system efficiency** Most
 336 of today’s FL efforts focus on either optimizing the statistical or
 337 the system efficiency, whereas we observe there exists a great
 338 need for jointly optimizing both efficiency: (1) practical FL suf-

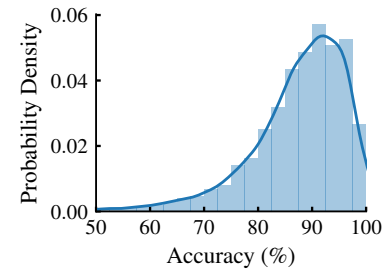


Figure 15: Biased accuracy distributions of the trained ShuffleNet model across clients.

339 fers biased model performance across clients (Figure 15). This can originate from the heterogeneous
340 data and system behaviors, because the system behavior determines the availability of client data
341 over training, wherein predicting this system behavior can curb the statistical drift in advance (e.g.,
342 prioritizing the use of upcoming offline clients). Moreover, the popular random client selection can
343 deemphasize clients with slow speed, leading to poor accuracy on slow clients; and (2) statistical
344 optimizations can leverage the heterogeneity nature of client system speed. For example, instead of
345 applying a one-fit-all strategy for all clients, faster workers can trade more system latency against
346 better statistical benefits. For example, faster workers can contribute larger but more accurate model
347 updates when using gradient compression.

348 6 Conclusion

349 To enable scalable, robust, and reproducible research of federated learning, we introduce FedScale,
350 a diverse set of realistic FL datasets in terms of scales, task categories and client system behaviors.
351 We provide realistic federated datasets for benchmarking today’s FL efforts. To enable efficient and
352 standardized FL evaluations, we introduce, FAR, a more scalable evaluation platform than the existing.
353 FAR performs fast-forward evaluation of the practical FL setting and produces FL runtime metrics
354 needed in today’s work. More subtly, FAR provides ready-to-use realistic datasets and flexible APIs
355 to allow more FL applications, such as benchmarking the performance of Neural Architecture Search,
356 model inference, and a broader view of federated data analytics (e.g., multi-party computation).

357 **Societal Impacts and Limitations** We expect FedScale to be a standard benchmark in federated
358 learning, contributing to the significant advancements of the field. One potential negative impact
359 is that FedScale might narrow down the scope of future papers to the tasks and dataset types that
360 have been included so far. In order to mitigate such a negative impact and limitation, we have made
361 FedScale open-source at: <https://github.com/SymbioticLab/FedScale>, and will regularly
362 update our datasets and tasks, based on the input from the community.

363 References

- 364 [1] AI Benchmark: All About Deep Learning on Smartphones. [http://ai-benchmark.com/
365 ranking_deeplearning_detailed.html](http://ai-benchmark.com/ranking_deeplearning_detailed.html).
- 366 [2] Common Voice Data. <https://commonvoice.mozilla.org/en/datasets>.
- 367 [3] Fox Go Dataset. <https://github.com/featurecat/go-dataset>.
- 368 [4] Google Open Images Dataset. [https://storage.googleapis.com/openimages/web/
369 index.html](https://storage.googleapis.com/openimages/web/index.html).
- 370 [5] iNaturalist 2019. [https://sites.google.com/view/fgvc6/competitions/
371 inaturalist-2019](https://sites.google.com/view/fgvc6/competitions/inaturalist-2019).
- 372 [6] LibriVox – Free public domain audiobooks. <https://librivox.org/>.
- 373 [7] MobiPerf. <https://www.measurementlab.net/tests/mobiperf/>.
- 374 [8] PySyft. <https://github.com/OpenMined/PySyft>.
- 375 [9] Reddit Comment Data. <https://files.pushshift.io/reddit/comments/>.
- 376 [10] Stack Overflow Data. [https://cloud.google.com/bigquery/public-data/
377 stackoverflow](https://cloud.google.com/bigquery/public-data/stackoverflow).
- 378 [11] Taobao Dataset. [https://tianchi.aliyun.com/dataset/dataDetail?dataId=56&
379 lang=en-us](https://tianchi.aliyun.com/dataset/dataDetail?dataId=56&lang=en-us).
- 380 [12] TensorFlow Federated. <https://www.tensorflow.org/federated>.
- 381 [13] Daniel J. Beutel, Taner Topal, Akhil Mathur, Xinchu Qiu, Titouan Parcollet, Pedro Porto Buarque
382 de Gusmao, and Nicholas D. Lane. FLOWER: A friendly federated learning framework. *arXiv
383 preprint arXiv:2007.14390*, 2021.
- 384 [14] Keith Bonawitz, Hubert Eichner, and et al. Towards federated learning at scale: System design.
385 In *MLSys*, 2019.
- 386 [15] Sebastian Caldas, Sai Meher, Karthik Duddu, and et al. Leaf: A benchmark for federated
387 settings. *NeurIPS’ Workshop*, 2019.
- 388 [16] Di Chai, Leye Wang, Kai Chen, and Qiang Yang. FedEval: A benchmark system with a
389 comprehensive evaluation model for federated learning. In *arxiv.org/abs/2011.09655*, 2020.
- 390 [17] Mingqing Chen, Rajiv Mathews, Tom Ouyang, and Françoise Beaufays. Federated learning of
391 out-of-vocabulary words. In *arxiv.org/abs/1903.10635*, 2019.

- 392 [18] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. EMNIST: an extension
393 of MNIST to handwritten letters. In *arxiv.org/abs/1702.05373*, 2017.
- 394 [19] Canh T. Dinh, Nguyen H. Tran, and Tuan Dung Nguyen. Personalized federated learning with
395 moreau envelopes. In *NeurIPS*, 2020.
- 396 [20] Hubert Eichner, Tomer Koren, H. Brendan McMahan, Nathan Srebro, and Kunal Talwar.
397 Semi-cyclic stochastic gradient descent. In *ICML*, 2019.
- 398 [21] Scott Ettinger, Shuyang Cheng, Benjamin Caine, Chenxi Liu, Hang Zhao, Sabeek Pradhan,
399 Yuning Chai, Benjamin Sapp, Charles Qi, Yin Zhou, Zoey Yang, Aurelien Chouard, Pei
400 Sun, Jiquan Ngiam, Vijay Vasudevan, Alexander McCauley, Jonathon Shlens, and Dragomir
401 Anguelov. Large scale interactive motion forecasting for autonomous driving : The waymo
402 open motion dataset. *CoRR*, abs/2104.10133, 2021.
- 403 [22] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning with
404 theoretical guarantees: A model-agnostic meta-learning approach. In *34th Conference on
405 Neural Information Processing Systems (NeurIPS 2020)*, 2020.
- 406 [23] David F. Fouhey, Weicheng Kuo, Alexei A. Efros, and Jitendra Malik. From lifestyle vlogs to
407 everyday interactions. In *CVPR*, 2018.
- 408 [24] Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. Inverting gradients
409 - how easy is it to break privacy in federated learning? In *NeurIPS*, 2020.
- 410 [25] Robin C. Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A
411 client level perspective. In *NeurIPS*, 2017.
- 412 [26] Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An efficient framework
413 for clustered federated learning. In *34th Conference on Neural Information Processing Systems
414 (NeurIPS 2020)*, 2020.
- 415 [27] Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An efficient framework
416 for clustered federated learning. In *NeurIPS*, 2020.
- 417 [28] Chaoyang He, Songze Li, Jinhyun So, and Xiao Zeng. FedML: A research library and bench-
418 mark for federated machine learning. In *arxiv.org/abs/2007.13518*, 2020.
- 419 [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
420 recognition. In *CVPR*, 2016.
- 421 [30] Sixu Hu, Yuan Li, Xu Liu, Qinbin Li, Zhaomin Wu, and Bingsheng He. The OARF
422 benchmark suite: Characterization and implications for federated learning systems. In
423 *arxiv.org/abs/2006.07856*, 2020.
- 424 [31] Peter Kairouz, Brendan McMahan, Shuang Song, Om Thakkar, Abhradeep Thakurta, and
425 Zheng Xu. Practical and private (deep) learning without sampling or shuffling. In
426 *arxiv.org/abs/2103.00039*, 2021.
- 427 [32] Peter Kairouz, H. Brendan McMahan, and et al. Advances and open problems in federated
428 learning. In *Foundations and Trends® in Machine Learning*, 2021.
- 429 [33] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U. Stich,
430 and Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for federated
431 learning. In *ICML*, 2020.
- 432 [34] Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian U Stich, and Martin Jaggi. Error feed-
433 back fixes signsgd and other gradient compression schemes. In *arXiv preprint arXiv:1901.09847*,
434 2019.
- 435 [35] Philipp Koehn. Europarl: A Parallel Corpus for Statistical Machine Translation. In *Conference
436 Proceedings: the tenth Machine Translation Summit*, 2005.
- 437 [36] Fan Lai, Xiangfeng Zhu, Harsha V. Madhyastha, and Mosharaf Chowdhury. Oort: Efficient
438 federated learning via guided participant selection. In *USENIX Symposium on Operating
439 Systems Design and Implementation (OSDI)*, 2021.
- 440 [37] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia
441 Smith. Federated optimization in heterogeneous networks. In *MLSys*, 2020.
- 442 [38] Wenqi Li, Fausto Milletari, and Daguang Xu. Privacy-preserving federated brain tumour
443 segmentation. In *Machine Learning in Medical Imaging*, 2019.
- 444 [39] Heiko Ludwig, Nathalie Baracaldo, Gegi Thomas, and et al. Ibm federated learning: An
445 enterprise framework white paper v0.1. In *arxiv.org/abs/2007.10987*, 2020.
- 446 [40] Peter Mattson, Christine Cheng, Cody Coleman, and et al. Mlperf training benchmark. In
447 *MLSys*, 2020.
- 448 [41] Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton van den Hengel. Image-based
449 recommendations on styles and substitutes. In *SIGIR*, 2015.

- 450 [42] H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas.
451 Communication-efficient learning of deep networks from decentralized data. In *AISTATS*, 2017.
- 452 [43] Matthias Paulik, Matt Seigel, Henry Mason, and et al. Federated evaluation and tuning for
453 on-device personalization: System design and applications. In *arxiv.org/abs/2102.08503*, 2021.
- 454 [44] Sashank Reddi, Zachary Charles, and et al. Adaptive federated optimization. In
455 *arxiv.org/abs/2003.00295*, 2020.
- 456 [45] Siva Reddy, Danqi Chen, and Christopher D. Manning. Coqa: A conversational question
457 answering challenge. *arXiv preprint arXiv:1808.07042*, 2019.
- 458 [46] Daniel Rothchild, Ashwinee Panda, Enayat Ullah, Nikita Ivkin, Ion Stoica, Vladimir Braverman,
459 Joseph Gonzalez, and Raman Arora. Fetchsgd: Communication-efficient federated learning
460 with sketching. In *ICML*, 2020.
- 461 [47] Mark Sandler, Andrew G. Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen.
462 Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, 2018.
- 463 [48] J. Schler, M. Koppel, S. Argamon, and J. Pennebaker. Effects of age and gender on blogging. In
464 *Proceedings of AAAI Spring Symposium on Computational Approaches for Analyzing Weblogs*,
465 2006.
- 466 [49] Gunnar A. Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav
467 Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In
468 *ECCV*, 2016.
- 469 [50] Ziteng Sun, Peter Kairouz, Ananda Theertha Suresh, and H Brendan McMahan. Can you really
470 backdoor federated learning. In *arxiv.org/abs/1911.07963*, 2019.
- 471 [51] Hongyi Wang, Kartik Sreenivasan, Shashank Rajput, Harit Vishwakarma, Saurabh Agarwal,
472 Jy yong Sohn, Kangwook Lee, and Dimitris Papailiopoulos. Attack of the tails: Yes, you really
473 can backdoor federated learning. In *NeurIPS*, 2020.
- 474 [52] Pete Warden. Speech commands: A dataset for limited-vocabulary speech recognition. In
475 *arxiv.org/abs/1804.03209*, 2018.
- 476 [53] Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2 a
477 large-scale benchmark for instance-level recognition and retrieval. In *arxiv.org/abs/2004.01804*,
478 2020.
- 479 [54] Chengxu Yang, Qipeng Wang, and et al. Characterizing impacts of heterogeneity in federated
480 learning upon large-scale smartphone data. In *WWW*, 2021.
- 481 [55] Timothy Yang, Galen Andrew, Hubert Eichner, Haicheng Sun, Wei Li, Nicholas Kong, Daniel
482 Ramage, and Françoise Beaufays. Applied federated learning: Improving Google keyboard
483 query suggestions. In *arxiv.org/abs/1812.02903*, 2018.
- 484 [56] Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J. Weiss, Ye Jia, Zhifeng Chen, and
485 Yonghui Wu. Libritts: A corpus derived from librispeech for text-to-speech. *arXiv preprint*
486 *arXiv:1904.02882*, 2019.
- 487 [57] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient
488 convolutional neural network for mobile devices. In *CVPR*, 2018.

489 **A Introduction of FedScale Datasets**

Category	Name	Data Type	#Clients	#Instances	Example Task
CV	iNature [5]	Image	2,295	193K	Classification
	FEMNIST [18]	Image	3,400	640K	Classification
	OpenImage [4]	Image	13,771	1.3M	Classification, Object detection
	Google Landmark [53]	Image	43,484	3.6M	Classification
	Charades [49]	Video	266	10K	Action recognition
	VLOG [23]	Video	4,900	9.6K	Classification, Object detection
	Waymo Motion [21]	Video	496,358	32.5M	Motion prediction
NLP	Europarl [35]	Text	27,835	1.2M	Text translation
	Blog Corpus [48]	Text	19,320	137M	Word prediction
	Stackoverflow [10]	Text	342,477	135M	Word prediction, Classification
	Reddit [9]	Text	1,660,820	351M	Word prediction
	Amazon Review [41]	Text	1,822,925	166M	Classification, Word prediction
	CoQA [45]	Text	7,189	114K	Question Answering
	LibriTTS [56]	Text	2,456	37K	Text to speech
	Google Speech [52]	Audio	2,618	105K	Speech recognition
Misc ML	Common Voice [2]	Audio	12,976	1.1M	Speech recognition
	Taobao [11]	Text	182,806	20.9M	Recommendation
	Fox Go [3]	Text	150,333	4.9M	Reinforcement learning

Table 4: Statistics of FedScale datasets. FedScale has 18 realistic client datasets, which are from the real-world collection, and we partitioned each dataset using its real client-data mapping.

490 FedScale currently has 18 realistic federated datasets across a wide range of scales and task categories
 491 (Table 4). Here, we provide the description of some representative datasets, and the reader can refer
 492 to FedScale repository (<https://github.com/SymbioticLab/FedScale>) for more datasets.

493 **Google Speech Commands.** A speech recognition dataset [52] with over ten thousand clips of
 494 one-second-long duration. Each clip contains one of the 35 common words (e.g., digits zero to nine,
 495 "Yes", "No", "Up", "Down") spoken by thousands of different people.

496 **OpenImage.** OpenImage [4] is a vision dataset collected from Flickr, an image and video hosting
 497 service. It contains a total of 16M bounding boxes for 600 object classes (e.g., Microwave oven). We
 498 clean up the dataset according to the provided indices of clients. In our evaluation, the size of each
 499 image is 256×256 .

500 **Reddit and StackOverflow.** Reddit [9] (StackOverflow [10]) consists of comments from the Reddit
 501 (StackOverflow) website. It has been widely used for language modeling tasks, and we consider each
 502 user as a client. In this dataset, we restrict to the 30k most frequently used words, and represent each
 503 sentence as a sequence of indices corresponding to these 30k frequently used words.

504 **VLOG.** VLOG [23] is a video dataset collected from YouTube. It contains more than 10k Lifestyle
 505 Vlogs, videos that people purportedly record to show their lives, from more than 4k actors. This
 506 dataset aimed at understanding everyday human interaction and contains labels for scene classification,
 507 hand-state prediction task, and hand detection.

508 **LibriTTS.** LibriTTS [56] is a large-scale text-to-speech dataset. It is derived from audiobooks that
 509 are part of the LibriVox project [6]. There are 585 hours of read English speech from 2456 speakers
 510 at 24kHz sampling rate.

511 **Taobao.** Taobao Dataset [11] is a dataset of click rate prediction about display Ad, which is
 512 displayed on the website of Taobao. It is composed of 1,140,000 users ad display/click logs for

513 8 days, which are randomly sampled from the website of Taobao. We partitioned it using its real
514 client-data mapping.

515 **Waymo Motion.** Waymo Motion [21] is composed of 103,354 segments each containing 20
516 seconds of object tracks at 10Hz and map data for the area covered by the segment. These segments
517 are further broken into 9 second scenarios (8 seconds of future data and 1 second of history) with 5
518 second overlap, and we consider each scenario as a client.

519 B Comparison with Existing FL Benchmarks

520 In this section, we compare FedScale with existing FL benchmarks in more details.

521 **Data Heterogeneity** Existing benchmarks for FL are mostly limited in the variety of realistic
522 datasets for real-world FL applications. Even they have various datasets (e.g., LEAF [15]) and
523 FedEval [16]), their datasets are mostly synthetically derived from conventional datasets (e.g.,
524 CIFAR) and limited to quite a few FL tasks. These statistical client datasets can not represent realistic
525 characteristics and are inefficient to benchmark various real FL applications. Instead, FedScale
526 provides 18 comprehensive realistic datasets for a wide variety of tasks and across small, medium,
527 and large scales, and these datasets can also be used in data analysis to motivate more FL designs.

528 **System Heterogeneity** The practical FL statistical performance also depends on the system het-
529 erogeneity (e.g., client system speed and availability of the client), which has inspired lots of
530 optimizations for FL system efficiency. However, existing FL benchmarks have largely overlooked
531 the system behaviors of FL clients, which can produce misleading evaluations, and discourages the
532 benchmarking of system efforts. To emulate the heterogeneous system behaviors in practical FL,
533 FedScale incorporates real-world traces of mobile devices, and associates each client with his system
534 speeds, as well as the availability. Moreover, it is non-trivial to emulate these behaviors at scale, so
535 we develop FAR, which is more efficient than the existing.

536 **Scalability** Existing frameworks, perhaps due to the heavy burden of building complicated sys-
537 tem support, largely rely on the traditional ML architectures (e.g., the primitive parameter-server
538 architecture of Pytorch). These architectures are primarily designed for the traditional large-batch
539 training on a number of workers, and each worker often trains a single batch at a time. However,
540 this is ill-suited to the simulation of thousands of clients concurrently: (1) they lack tailored system
541 implementations to orchestrate the synchronization and resource scheduling, for which they can easily
542 run into synchronization/memory issues and crash down; (2) their resource can be under-utilized, as
543 FL evaluations often use a much smaller batch size than that in the traditional architecture.

544 Tackling all these inefficiencies requires domain-specific system designs, and the FAR is refactored
545 atop of our Oort project [36]. Specifically, we first built an advanced resource scheduler: It monitors
546 the fine-grained resource utilization of machines, queues the overcommitted simulation requests,
547 adaptively dispatches simulation requests of the client across machines to achieve load balance,
548 and then orchestrates the simulation based on the client mirror clock. Moreover, given a much
549 smaller batch size in FL, we maximize the resource utilization by overlapping the communication
550 and computation phrases of different client simulations. The former and the latter make FedScale
551 more scalable across machines and on single machines, respectively.

552 **Modularity** As shown in Table 1, some existing frameworks (e.g., LEAF and FedEval) do not
553 provide user-friendly modularity, which requires great engineering efforts to benchmark different
554 components, and we recognize that FedML and Flower provide the API modularity in this table too.

555 On the other hand, FAR’s modularity for easy deployments and broader use cases is not limited to
556 APIs (Figure 7): (1) FAR Data Loader: it simplifies and expands the use of realistic datasets. e.g.,
557 developers can load and analyze the realistic FL data to motivate new algorithm designs, or imports
558 new datasets/customize data distributions in FedScale evaluations; (2) Client simulator: it emulates
559 the system behaviors of FL clients, and developers can customize their system traces in evaluating
560 the FL system efficiency too; (3) Resource Manager: it hides the system complexity in training
561 large-scale participants simultaneously for the deployment.

```

from fedscale.core.client_manager import ClientManager
import Oort

class Customized_ClientManager(ClientManager):
    def __init__(self, *args):
        super().__init__(*args)
        self.oort_selector = Oort.create_training_selector(*args)

    # Replace default client selection algorithm w/ Oort
    def resampleClients(self, numofClients, cur_time, feedbacks):
        # Feed Oort w/ execution feedbacks from last training round
        oort_selector.update_client_info(feedbacks)
        selected_clients = oort_selector.select_participants(numofClients, cur_time)

    return selected_clients

```

Figure 16: Evaluate new client selection algorithm [36].

<hr/> <pre> from fedscale.core.client import Client class Customized_Client(Client): # Customize the training on each client def train(self, client_data, model, conf): # Get the training result from # the default training component training_result = super().train(client_data, model, conf) # Implementation of compression compressed_result = compress_impl(training_result) return compressed_result </pre> <hr/>	<hr/> <pre> from fedscale.core.client import Client class Customized_Client(Client): # Customize the training on each client def train(self, client_data, model, conf): # Get the training result from # the default training component training_result = super().train(client_data, model, conf) # Clip updates and add noise secure_result = secure_impl(training_result) return secure_result </pre> <hr/>
--	---

Figure 17: Evaluate model compression [46]. Figure 18: Evaluate security enhancement [50].

562 C Examples of New Plugins

563 In this section, we demonstrate several examples to show the ease of integrating today’s FL efforts
564 for realistic evaluations in FedScale.

565 At its core, FAR provides flexible APIs on each module so that the developer can access and customize
566 methods of the base class. Note that FAR will automatically integrate new plugins into evaluations,
567 and then produces practical FL metrics. Figure 16 demonstrates that we can easily evaluate new client
568 selection algorithms, Oort [36], by modifying a few lines of the `clientManager` module. Similarly,
569 Figure 17 and Figure 18 show that we can extend the basic `Client` module to apply new gradient
570 compression [46] and enhancement for malicious attack [50], respectively.

571 D FedScale Maintenance

572 **Availability of data and platform** We have made FedScale open-source on the Github (<https://github.com/SymbioticLab/FedScale>). So the code and dataset can be downloaded from
573 this repository. For each dataset, we provide detailed descriptions (README.md) of the source,
574 organization, format and use case under the repository. So far, these datasets are host on Dropbox,
575 and we are migrating them to the stable storage of AWS. For the evaluation platform FAR, we provide
576 the configuration and job submission guideline as well. We encourage the reviewer to refer to our
577 repository for more details.
578

579 **Maintenance plan and responsibility.** We are actively updating our benchmark weekly, based on
580 the feedback from the community. Currently, our dataset and platform are subject to the *Apache-2.0*

581 *License.* We respect the contributor of each dataset in following ways: (1) we provide the scripts for
582 the developer to preprocess the downloaded raw data from its original source. This will absolutely
583 obey the rule of each contributor; (2) for those publicly available and widely-used dataset, we
584 temporarily host the processed data on our repository. However, we are creating permissive license
585 for each dataset and acknowledgment to respect their contributor, and highlight all assets in our
586 repository are for research purpose only; (3) for all assets in our work, we have removed the sensitive
587 information and use anonymous information to partition the data; (4) we are keeping in touch with all
588 the contributors, and will fix any issues (e.g., by removing that dataset) once that happens.