

445 A Co-Adaptation in Subset-based FL

446 Federated Dropout is originally inspired by regular dropout [40], a regularization technique that
 447 constrains the capacity of a large NN model by randomly dropping parameters from the training,
 448 thereby limiting co-adaptation among parameters of the model. This is essential to improve the
 449 accuracy and reduce the over-fitting of parameters, as shown in various studies. FD and FedRolex
 450 adopt the dropout technique, removing CNN’s filters in a round-based manner. These techniques,
 451 however, exercise dropout to its extreme, dropping a large part of filters so that not enough co-
 452 adaptation between filters remains. In particular, the gradients for the subset of parameters trained on
 453 a device are calculated without consideration of the error of the remaining parameters that reside on
 454 the server. These subsets are randomly changing over time and devices, reducing the co-adaptation of
 455 this distributed training process. Add to these the fact that the data is also distributed over devices,
 456 so applying such a random scheme significantly decreases the chance that a subset of parameters is
 457 being trained together over a sizable proportion of the data.

458 To further study the effects of co-adaptation on the reachable accuracy and the differences between FD
 459 and FedRolex, we run the following experiment, using CIFAR10 with ResNet20 and $s_{\text{FD}/\text{FedRolex}} =$
 460 0.25:

- 461 • We modify FD s.t. all devices train the same random subset per round, i.e., the same
- 462 indices $I^{(r,c)} = I^{(r)}$ per round (index k is omitted for simplicity).
- 463 • We limit the randomness by randomly selecting $I^{(r)}$ out of a set $\mathcal{I} = \{I_1, \dots, I_{|\mathcal{I}|}\}$ of
- 464 randomly initialized subsets that are generated once prior to training. Each round, $I^{(r)}$ is a
- 465 random selection out of \mathcal{I} .

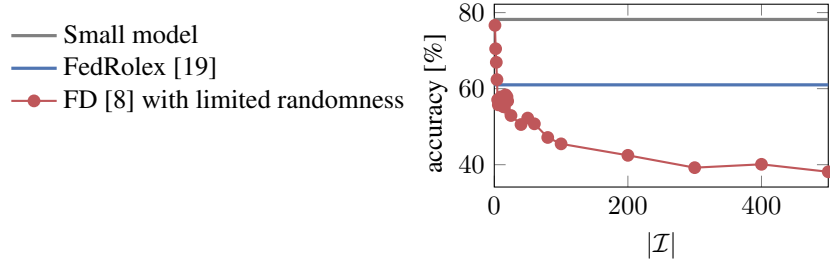


Figure 5: FD with limited randomness using CIFAR10 with ResNet20 and $s_{\text{FD}/\text{FedRolex}} = 0.25$.

466 We vary the randomness by varying the number of sampled subsets in \mathcal{I} , s.t. $|\mathcal{I}| \in [0, 500]$. Thereby,
 467 the probability of a specific subset being selected is $p = \frac{1}{|\mathcal{I}|}$. Therefore, all the devices train the same
 468 submodel in a round, and there are only $|\mathcal{I}|$ submodels that would be trained by devices over the
 469 training period. Evaluation is done with the full server model.

470 We observe the following effects: Firstly, the final accuracy drops proportionally to $\sim p$. Secondly,
 471 it can be observed that in the case of $|\mathcal{I}| = 1$, FD behaves similarly to the small model baseline,
 472 as always the same subset is used for training. We also observe that remaining untrained filters
 473 have a minor effect on the accuracy when compared with a small model. However, because of these
 474 untrained parameters, the model fails to reach higher accuracies as with SLT (see Section 3). The
 475 accuracy drops with introducing more randomness to the training process (i.e., increasing $|\mathcal{I}|$). Lastly,
 476 it can be observed that the rolling window approach of FedRolex is a special case of FD with limited
 477 randomness (i.e., $|\mathcal{I}| = 5$ in this experiment).

478 B Training Memory Measurements in PyTorch

479 We measure the maximum memory requirement for the evaluated NN models ResNet and DenseNet
 480 using PyTorch 1.10. Specifically, we measure the size of the activations, gradients, and weights.
 481 These memory measurements are done offline (prior to training) and do not require any data.

- 482 • **Measurement of weights:** To measure the size of the weights, we sum up all tensors that
- 483 are present in the NN’s `state_dict`.

484 • **Measurement of activations and gradients:** To measure the size of the activations that have
485 to be kept in memory, as well as the gradients, we apply `backward_hooks` to all relevant
486 PyTorch modules in an NN. Specifically, we add these hooks to `Conv2d`, `BatchNorm2d`,
487 `ReLU`, `Linear`, and `Add` operations. If a hook attached to a module is called, we add the
488 respective size of the activation map and the size of the calculated gradient to a global
489 variable add up all required activations and gradients.