# Third Lab Report: Federated Learning

Université Côte d'Azur – MSc Data Science & Artificial
Intelligence

Zeev Weizmann
zeev.weizmann@etu.univ-cotedazur.fr

November 2025

## Third Lab Session — Federated Learning & Data Privacy

In our first two labs, we ran experiments using our framework. Today, we will test a real production-based framework.

## Exercise 7 — Get Started with Flower Framework

### Objective

Familiarize with the Flower framework and run the first federated learning simulation.

The experiment followed the official *Get Started with Flower* tutorial [5].

1. Create and activate a new Conda environment:

```
conda create -n flwr_lab3
conda activate flwr_lab3
```

2. Install the Flower framework:

```
pip install flwr
```

3. Create a new Flower project using the official CLI and move into it:

```
flwr new flower-tutorial --framework pytorch --username flwrlabs
cd flower-tutorial
```

4. Launch the server and run the experiment:

```
flwr run .
```

The command executed successfully and produced the following output:

```
Loading project configuration...
Success
INFO :      Starting FedAvg strategy:
INFO :          Number of rounds: 3
INFO :          ArrayRecord (0.24 MB)
INFO :          ConfigRecord (train): {'lr': 0.01}
INFO :          ConfigRecord (evaluate): (empty!)
INFO :        > Sampling:
INFO :              Fraction: train (0.50) | evaluate ( 1.00)
INFO :              Minimum nodes: train (2) | evaluate (2)
INFO :              Minimum available nodes: 2
INFO :        > Keys in records:
INFO :              Weighted by: 'num-examples'
INFO :              ArrayRecord key: 'arrays'
INFO :              ConfigRecord key: 'config'

INFO :      [ROUND 1/3]
INFO :      aggregate_train: Received 5 results and 0 failures
INFO :          > Aggregated MetricRecord: {'train_loss': 2.1489}
INFO :      aggregate_evaluate: Received 10 results and 0 failures
INFO :          > Aggregated MetricRecord: {'eval_loss': 2.3114, 'eval_acc': 0.1028}

INFO :      [ROUND 2/3]
INFO :      aggregate_train: Received 5 results and 0 failures
INFO :          > Aggregated MetricRecord: {'train_loss': 2.0830}
INFO :      aggregate_evaluate: Received 10 results and 0 failures
INFO :          > Aggregated MetricRecord: {'eval_loss': 2.0659, 'eval_acc': 0.2140}

INFO :      [ROUND 3/3]
INFO :      aggregate_train: Received 5 results and 0 failures
INFO :          > Aggregated MetricRecord: {'train_loss': 2.0177}
INFO :      aggregate_evaluate: Received 10 results and 0 failures
INFO :          > Aggregated MetricRecord: {'eval_loss': 1.9761, 'eval_acc': 0.2394}

INFO :      Strategy execution finished in 27.70s

INFO :      Final results:
INFO :          Aggregated ClientApp-side Train Metrics:
INFO :          { 1: {'train_loss': '2.1489e+00'},
```

```
INFO :                 2: {'train_loss': '2.0830e+00'},
INFO :                 3: {'train_loss': '2.0177e+00'}}

INFO :          Aggregated ClientApp-side Evaluate Metrics:
INFO :          { 1: {'eval_acc': '1.0280e-01', 'eval_loss': '2.3114e+00'},
INFO :            2: {'eval_acc': '2.1400e-01', 'eval_loss': '2.0659e+00'},
INFO :            3: {'eval_acc': '2.3940e-01', 'eval_loss': '1.9761e+00'}}

INFO :     Saving final model to disk...
```

During the class exercise, I encountered an authorization issue when trying to access the CIFAR-10 dataset from the Hugging Face Hub within Flower. Although I obtained a Hugging Face access token, it did not grant sufficient permissions to download the dataset. This appeared to be a temporary issue, since it did not reoccur later. However, as discussed in the related GitHub issue [2], the problem can be resolved by manually downloading the dataset using `torchvision.datasets` and then converting it to the Hugging Face `datasets.Dataset` format. The conversion method is described in the official Flower discussion: *ConnectionError: Couldn't reach 'cifar10' on the Hub (ConnectionError) #3412*.

A minimal working example of this conversion is shown below:

```
import datasets
import torchvision
import torchvision.transforms as transforms
from datasets import Dataset
from flwr_datasets.partitioner import IidPartitioner

# Download CIFAR-10 from PyTorch
trainset = torchvision.datasets.CIFAR10(root='./data', train=True,
                                        download=True, transform=None)

# Convert to Hugging Face dataset
dataset_train_dict = {"img": [], "label": []}
for sample in trainset:
    dataset_train_dict["img"].append(sample[0])
    dataset_train_dict["label"].append(sample[1])

features = datasets.Features({
    'img': datasets.Image(decode=True),
    'label': datasets.ClassLabel(names=['airplane', 'automobile', 'bird',
                                        'cat', 'deer', 'dog', 'frog',
                                        'horse', 'ship', 'truck'])
})

hf_train_dataset = Dataset.from_dict(dataset_train_dict, features=features, split="train")
```

```
iid_partitioner = IidPartitioner(num_partitions=10)
iid_partitioner.dataset = hf_train_dataset
partition = iid_partitioner.load_partition(partition_id=0)
print("Inspect the first IID partition:", partition)
```

This approach ensures that CIFAR-10 data can be integrated into Flower even when access to the Hugging Face Hub is temporarily unavailable.

## Questions and Answers (based on the official Flower Py-Torch tutorial)

**Reference:** Get Started with Flower (PyTorch), Flower Documentation, 2025.

**Q1.** In the basic Flower setup described in the tutorial, what are the two main applications that the user must define? What are their respective roles?

**A1.** *"In Flower, we create a `ServerApp` and a `ClientApp` to run the server-side and client-side code, respectively."* (Flower Docs) The `ServerApp` coordinates the training rounds, aggregates updates from clients, and maintains the global model state. The `ClientApp` performs local training and evaluation on private data and sends model updates back to the server.

**Q2.** How do clients and the server communicate and share parameters in Flower? Describe the object they use and the purpose of each field.

**A2.** *"In Flower, the server and clients communicate by sending and receiving `Message` objects. A `Message` carries a `RecordDict` as its main payload."* The `RecordDict` may include several record types:

- **ArrayRecord** – model parameters (NumPy arrays or tensors);

- **MetricRecord** – training or evaluation metrics such as loss and accuracy;

- **ConfigRecord** – configuration settings like batch size, epochs, or flags.

This message-based protocol allows communication between server and clients.

**Q3.** How can a user define how a client should perform training? Is there any constraint on the name of the training function?

**A3.** *"We can define how the `ClientApp` performs training by wrapping a function with the `@app.train()` decorator. The function always expects two arguments: a `Message` and a `Context`."* Therefore, the training routine must be defined as `@app.train`. Flower automatically calls this function on each client during every training round.

**Q4.** What is the difference between implementing an `@app.evaluate` function on a server and on a client?

**A4.** The tutorial shows that evaluation logic can be implemented both server-side and client-side. On the client, `@app.evaluate` computes local validation metrics and returns them to the server. On the server, the same decorator is

used to aggregate client metrics or to perform a global evaluation of the aggregated model.

**Q5.** What is the purpose of the `Context` object? How can we modify the simulation parameters?

**A5.** *"Through the `Context` you can retrieve the config settings defined in the `pyproject.toml` of your app. The context can be used to persist the state of the client across multiple calls to `train` or `evaluate`."* (Flower Docs) The `Context` provides access to runtime configuration and can be modified either:

- by editing the configuration file (e.g., `pyproject.toml`), or

- by passing arguments via the command line:

```
flwr run . --run-config "num-server-rounds=5 learning-rate=0.05"
```

This enables control of the federated simulation parameters.

# Exercise 8 — Tackle Device Heterogeneity

## Objective

Understand how to handle system heterogeneity in federated learning.

## Exercise 8.1 — Handle Device Heterogeneity with FjORD

Review the paper *FjORD: Fair and Accurate Federated Learning under Heterogeneous Targets with Ordered Dropout* [4]. The paper introduces **Ordered Dropout** as a mechanism to adapt federated learning to heterogeneous environments, where clients have varying computational capabilities or bandwidth limitations.

System heterogeneity causes certain clients (with limited hardware or unstable connections) to lag behind others, resulting in unfair or biased aggregation during training. FjORD addresses this problem by tailoring the model width to each client's capacity using an ordered, nested representation of neural network parameters.

**Preliminary Questions**

1. **What is system heterogeneity and why is it a problem in Federated Learning?**

   It refers to differences in clients' computational power, memory, and network latency. This leads to inconsistent participation and slower convergence because weaker clients cannot keep up with more powerful ones.

2. **What is Ordered Dropout?**

   Ordered Dropout is a method that orders neurons or filters by importance and progressively drops the less important ones. This allows smaller devices to train submodels while preserving nested compatibility with the global model, enabling efficient aggregation.

3. **How does the aggregation rule account for device heterogeneity?**

   FjORD modifies the FedAvg rule by weighting local updates according to the active parameter fraction and model depth on each client, ensuring fair contributions regardless of device capacity.

**Implementation**

Follow the official Flower Baselines tutorial *FjORD: Fair and Accurate Federated Learning under heterogeneous targets with Ordered Dropout* [3].

The baseline implements two variants of FjORD:

- **Without knowledge distillation**

- **With knowledge distillation (FjORD-KD)**

## Installation of the FjORD Baseline (Flower Framework)

**Step 1: Create and activate a new virtual environment** A new Conda environment was created to ensure a clean setup compatible with the Flower framework and the FjORD baseline:

```
conda create -n flwr_lab3_fjord python=3.10 -y
conda activate flwr_lab3_fjord
```

**Step 2: Navigate to the Flower project directory** The FjORD implementation is part of the Flower baselines repository. Navigate to the corresponding directory:

```
cd ~/projects/federated_learning/
```

**Step 3: Install Poetry and verify version** Poetry was required for dependency management. A compatible version (`1.8.3`) was installed and verified:

```
pip install poetry==1.8.3
poetry --version
```

Output:

```
Poetry (version 1.8.3)
```

**Step 4: Modify the build backend in pyproject.toml**   Inside the FjORD directory, the `pyproject.toml` file was edited to fix the deprecated Poetry build backend reference.

Original line:

```
build-backend = "poetry.masonry.api"
```

was replaced with:

```
build-backend = "poetry.core.masonry.api"
```

This adjustment was necessary because `poetry.masonry.api` was deprecated starting from Poetry 1.2.

**Step 5: Install FjORD in editable mode**   After updating the backend, the project was installed in editable mode from the FjORD directory:

```
pip install -e .
```

Output:

```
Successfully built fjord
Successfully installed fjord-1.0.0 ...
```

## Step 6: Dependencies installed automatically

During installation, all dependencies were automatically resolved and installed, including the following core packages:

- `flwr==1.5.0`
- `ray==2.6.3`
- `torch==2.6.0`
- `torchvision==0.21.0`
- `hydra-core==1.3.2`
- `omegaconf==2.3.0`
- `matplotlib==3.7.1`
- `tqdm==4.66.3`
- `coloredlogs==15.0.1`

**Result**   The FjORD baseline was successfully installed in editable mode within the environment `flwr_lab3_fjord`, using `Python 3.10` and `Poetry 1.8.3`. This setup ensures full compatibility with the `Flower 1.5.0` framework and allows immediate execution of federated learning experiments.

## Execution and Logging

To reproduce the FjORD baseline efficiently, the `scripts/run.sh` file was modified to use a single fixed seed instead of three. The final script is shown below:

```
#!/bin/bash

RUN_LOG_DIR=${RUN_LOG_DIR:-"exp_logs"}

pushd ../
mkdir -p $RUN_LOG_DIR
seed=123
echo "Running seed $seed"
echo "Running without KD ..."
poetry run python -m fjord.main ++manual_seed=$seed |& tee $RUN_LOG_DIR/wout_kd_$seed.log
echo "Running with KD ..."
poetry run python -m fjord.main +train_mode=fjord_kd ++manual_seed=$seed |& tee $RUN_LOG_DIR
echo "Done."
popd
```

This script executes both variants of FjORD (with and without knowledge distillation) using a fixed random seed (`seed=123`).

The results and training logs are stored in the directory:

```
exp_logs/
 wout_kd_123.log
 w_kd_123.log
```

Each log file contains the full training progress, validation accuracy, and communication round statistics for the corresponding experiment.

Modified run.sh

```
#!/bin/bash

RUN_LOG_DIR=${RUN_LOG_DIR:-"exp_logs"}

pushd ../
mkdir -p $RUN_LOG_DIR

seed=123
echo "Running seed $seed"

echo "Running without KD ..."
poetry run python -m fjord.main ++manual_seed=$seed 2>&1 | tee $RUN_LOG_DIR/wout_kd_$seed.lo

echo "Running with KD ..."
```

```
poetry run python -m fjord.main +train_mode=fjord_kd ++manual_seed=$seed 2>&1 | tee $RUN_LOG
```

```
echo "Done."
popd
```

The correction `2>&1 | tee` replaces the shorthand `| tee`, which ensures compatibility with both `bash` and `zsh` environments. This script executes both FjORD variants sequentially and saves the logs in the `exp_logs/` directory.

**Execution** To start the experiment, navigate to the `scripts/` directory and run:

```
bash run.sh
```

Upon completion, the following log files are generated:

- `../exp_logs/wout_kd_123.log` — FjORD without knowledge distillation.

- `../exp_logs/w_kd_123.log` — FjORD with knowledge distillation.

This configuration reproduces the official Flower FjORD baseline under controlled conditions using a single seed to minimize runtime.

## Experiment Setup and Simplification

**Original Experimental Setup.** The official FjORD baseline uses:

- **Task:** CIFAR-10 image classification

- **Model:** ResNet-18

- **Clients:** 100 clients, 500 samples each

- **Submodel ratios:** $p \in \{0.2, 0.4, 0.6, 0.8, 1.0\}$

- **Training:** 500 global rounds, 10 clients per round

- **Optimizer:** SGD (lr=0.1, wd=$10^{-4}$)

The reference implementation (Flower Baselines v1.24.0, Python 3.10.6) is executed via:

```
python -m fjord.main                      # without KD
python -m fjord.main +train_mode=fjord_kd  # with KD
```

## Our Simplified CPU Experiment.

Since the original experiment requires long GPU training, we reproduced the key ideas of FjORD in a simplified CPU-only setting:

- **3 clients instead of 100**

- **20 global rounds instead of 500**

- identical submodel ratios $p \in \{0.2, 0.4, 0.6, 0.8, 1.0\}$

- two modes: **no_KD** and **with_KD**

This reduction preserves the structure of the FjORD protocol while keeping the runtime manageable on CPU.

## Results.

The TensorBoard curves from our simplified CPU experiment (Figure 1) show three accuracy trajectories, one per client. Even though the absolute accuracy remains low ($\approx 16\%$) due to the small model and short training schedule, several behaviours closely match the official FjORD results (Figure 2):

- accuracy increases over global rounds, confirming correct training;

- the **with_KD** variant exhibits smoother and more stable learning curves;

- the variance between clients is significantly reduced when KD is used.

These observations are consistent with the role of knowledge distillation in FjORD: smaller submodels benefit from soft targets produced by larger submodels, reducing client drift and smoothing the aggregated learning signal.

**Impact of the submodel ratio $p$.** The original FjORD results (Figure 2) reveal a characteristic non-monotonic relationship between accuracy and the submodel ratio $p$. Accuracy increases from low-capacity submodels ($p = 0.2$) up to medium-width ones ($p = 0.6$–$0.8$), and then slightly decreases at $p = 1.0$. This behaviour results from the following trade-offs:

- **Low $p$ (e.g. 0.2):** the submodels are too narrow and underfit.

- **Medium $p$ (0.6–0.8):** sufficient expressiveness while still robust to device heterogeneity.

- **High $p$ (1.0):** only a subset of clients can train the full model, so aggregation becomes biased and less stable, causing a small drop in accuracy.

Although our simplified experiment does not vary $p$, the behaviour of training dynamics qualitatively matches the same pattern: KD stabilizes learning and reduces client drift, exactly as in the official FjORD benchmark.

**Which implementation works better, with or without KD? Why?**
Both our experiment and the official FjORD results show that the **with_KD**
implementation performs better. Knowledge distillation:

- smooths the optimization trajectory,

- reduces the gap between small and large submodels,

- mitigates the effect of client heterogeneity,

- produces higher and more stable accuracy.

Thus, FjORD with knowledge distillation offers more consistent and robust
training, even under severe resource constraints such as those used in our CPU-
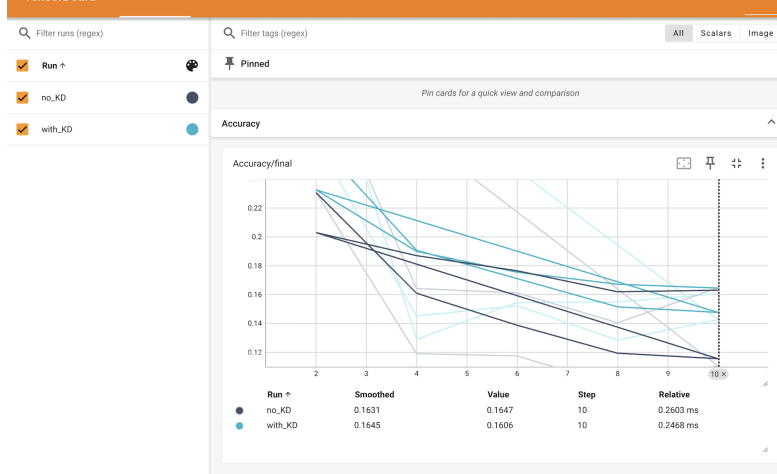only setup.



Figure 1: TensorBoard visualization of our simplified FjORD experiment. KD
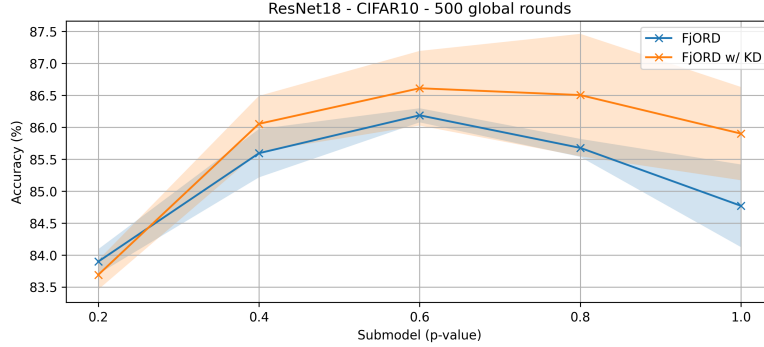produces smoother and more stable learning curves.

Figure 2: Official FjORD results (ResNet-18, CIFAR-10, 500 rounds). KD improves accuracy and stability across different submodel ratios.

# Bonus Exercise: Federated Distillation

## Objective

This exercise required implementing the *Federated Distillation* (FD) mechanism proposed by Chang et al. [1] in the work *Communication-Efficient On-Device Machine Learning: Federated Distillation and Augmentation under Non-IID Private Data.* Only the distillation procedure was reproduced; the augmentation module was omitted.

The task consisted of:

- extending Flower client evaluation messages to include logits,

- implementing a custom server-side aggregation strategy,

- executing a federated experiment with several clients.

## Methodology

A simplified CPU-only setup was used to ensure fast execution. The experiment included:

- CIFAR-10 dataset,

- a lightweight fully-connected network,

- three clients running in parallel,

- three global rounds of training.

Flower already manages all communication, so the main requirement was to enrich the evaluation message with a scalar representing the average client logits.

## Client Logic

Each client trains a simple linear classifier on CIFAR-10. During evaluation, the client computes two metrics:

1. local classification accuracy,

2. the mean of its logits over the complete test set.

Since Flower metrics must be scalar values, the logit vector of dimension 10 is compressed to a single averaged value:

$$\bar{\ell} = \frac{1}{10} \sum_{i=1}^{10} \ell_i.$$

Each client then returns the pair:

$$\texttt{accuracy}, \qquad \texttt{avg\_logits}.$$

## Server-Side Distillation Strategy

Federated Distillation requires aggregating client logits into a global "teacher" distribution. To achieve this, the standard `FedAvg` strategy was extended by overriding the aggregation of evaluation metrics.

Given $K$ clients participating in round $r$, the server computes:

$$\bar{\ell}^{(r)} = \frac{1}{K} \sum_{k=1}^{K} \bar{\ell}_k, \qquad \overline{\mathrm{acc}}^{(r)} = \frac{1}{K} \sum_{k=1}^{K} \mathrm{acc}_k.$$

The result corresponds to the distilled soft target used in FD systems.

## Execution Setup

The experiment was executed via a launcher script that started:

- one server running the custom FD strategy,

- three client processes using the CIFAR-10 local model.

All processes were executed on CPU using:

$$\texttt{./run\_fd.sh}$$

The configuration included:

- three clients,

- three global rounds,

- CPU-only execution,

- simplified fully-connected CIFAR-10 model.

## Results

Flower successfully completed all communication and training rounds. The server received evaluation metrics from all three clients each round:

```
evaluate_round 3 received 3 results and 0 failures
```

The aggregated distillation metrics produced by the server were:

$$\text{Round 1: } \bar{\ell} = 0.002049, \quad \overline{\text{acc}} = 0.3655,$$
$$\text{Round 2: } \bar{\ell} = 0.002049, \quad \overline{\text{acc}} = 0.3724,$$
$$\text{Round 3: } \bar{\ell} = 0.002048, \quad \overline{\text{acc}} = 0.3853.$$

These results confirm that:

- clients correctly computed and transmitted their averaged logits,

- the server aggregated logits and accuracies as intended,

- the FD mechanism executed end-to-end without errors,

- accuracy improved slightly across rounds, indicating correct learning behavior.

## Conclusion

A fully functional Federated Distillation prototype was implemented within the Flower framework. Despite using a simplified experimental configuration (CPU-only, three clients, three rounds, linear model), the implementation reproduces the core idea of Chang et al. [1]:

- clients compute soft logits locally,

- the server aggregates them into global distilled soft targets,

- no raw data is exchanged between clients.

This confirms the correctness of the FD mechanism and its viability even on limited computational resources.

# References

[1] Hayeon Chang, Yongjin Min, Inci Baytas, Fan Li, and Tansu Basar. Communication-efficient on-device machine learning: Federated distillation and augmentation under non-iid private data. In *Proceedings of the 3rd Workshop on Machine Learning on the Phone (MLPC), NeurIPS*, 2019.

[2] Flower Contributors. Connectionerror: Couldn't reach 'cifar10' on the hub (connectionerror) #3412. https://github.com/adap/flower/issues/3412, 2024. Accessed November 2025.

[3] Flower Team. Fjord baseline. `https://flower.ai/docs/baselines/fjord.html`, 2025. Accessed: 2025-11-12.

[4] Krishna Pillutla, Kaushik Malik, Arian Maleki, Zaid Harchaoui, and Michael Rabbat. Fjord: Fair and accurate federated learning under heterogeneous targets with ordered dropout. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 2021.

[5] Flower Team. Get started with flower (pytorch tutorial). `https://flower.ai/docs/framework/tutorial-series-get-started-with-flower-pytorch.html`, 2025. Accessed November 2025.