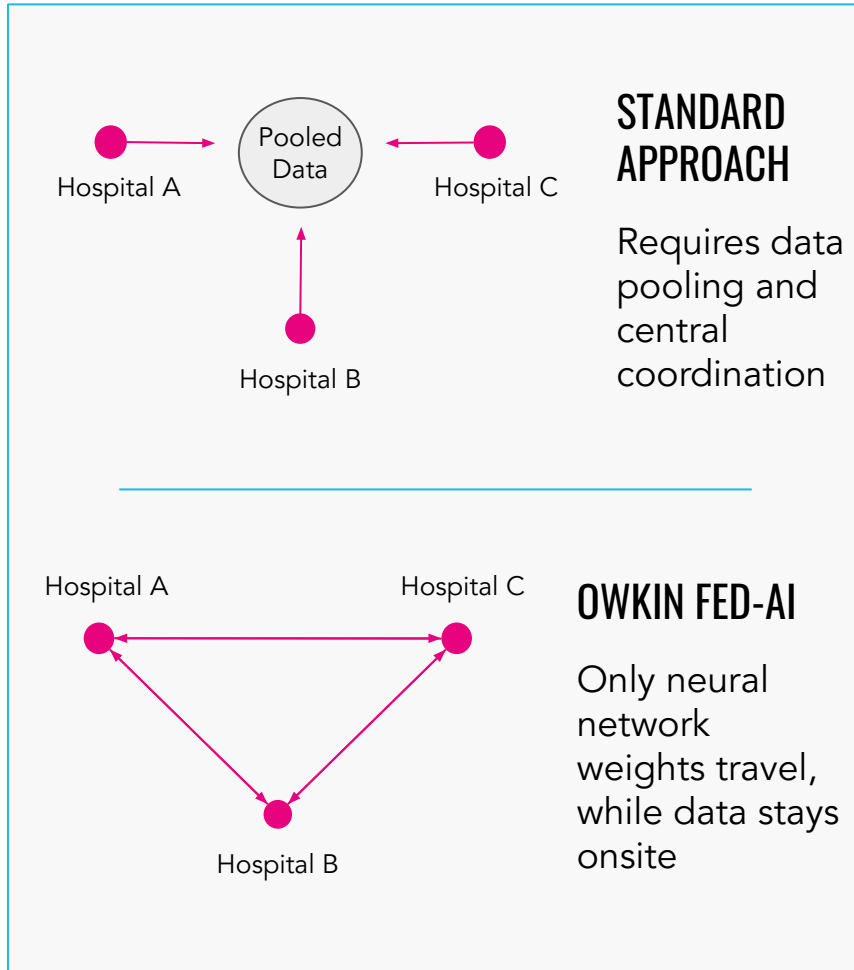# Distributed learning on sensitive health data

2020/01/27

Camille Marini

OWKIN

# FEDERATED LEARNING

Owkin develops federated learning to train machine learning models on distributed data at scale



**STANDARD APPROACH**

Requires data pooling and central coordination

**OWKIN FED-AI**

Only neural network weights travel, while data stays onsite

Data pooling implies a reduced control and governance of data owners

No transparency on how the algorithm is trained and how the data are used

# OUR PRODUCT: OWKIN CONNECT

Owkin Connect helps build state-of-the-art models from heterogeneous multicentric data, with high privacy and traceability standards.

Owkin Connect connects data managers to our data scientists

| Privacy | Traceability |
| --- | --- |

# Real World Deployments

# FED-AI REAL WORLD DEPLOYMENTS

## Pioneering FED-AI to connect distributed data at scale

- Won €10M BPI-funded Healthchain grant

- Consortium of 7 leading academic medical centers, 2 research centers and 2 start-ups
- Initial collaborations to develop predictive models of treatment outcomes in Breast Cancer & Melanoma
- Project coordination done by Owkin

- Live working experiment on real data

*« This project is supported by Bpifrance as part of the "Healthchain" project, which resulted from the "Digital Investments Program for the major challenges of the future" RFP. As part of the "Healthchain" project, a consortium coordinated by Owkin (a private company) has been established, including the Substra association, Apricity (a private company), the Assistance Publique des Hôpitaux de Paris, the University Hospital Center of Nantes, the Léon Bérard Center, the French National Center for Scientific Research, the École Polytechnique, the Institut Curie and the University of Paris Descartes ».*

# FED-AI REAL WORLD DEPLOYMENTS

Pioneering FED-AI to connect distributed data at scale

Project website

# Technical description

# OWKIN CONNECT - OPEN CORE APPROACH

Powered by the Open-Source Framework Substra

PROPRIETARY

Interfaces for Data Scientists

ML algorithms

Data preparation tools

Federated Learning strategies

Privacy preserving methods

**OPEN SOURCE**
Setup of network nodes
Trustless traceability and orchestration
Execution of ML tasks
Simple programmatic interfaces

SUBSTRA

GitHub SubstraFoundation

# OWKIN CONNECT

OWKIN
# OWKIN CONNECT

**Framework for ML orchestration
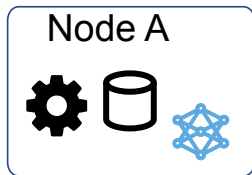on decentralized sensitive data**

**Data privacy**   **Traceability**

Data
agnostic          Algorithm
agnostic          ML framework agnostic

OWKIN
# ASSETS



**Node A**

**Node B**

**Node C**

**Node D**

- Objective
    - scientific question
    - evaluation metrics
    - test dataset

- Dataset
    - set of data samples
    - opener (script to read data samples)
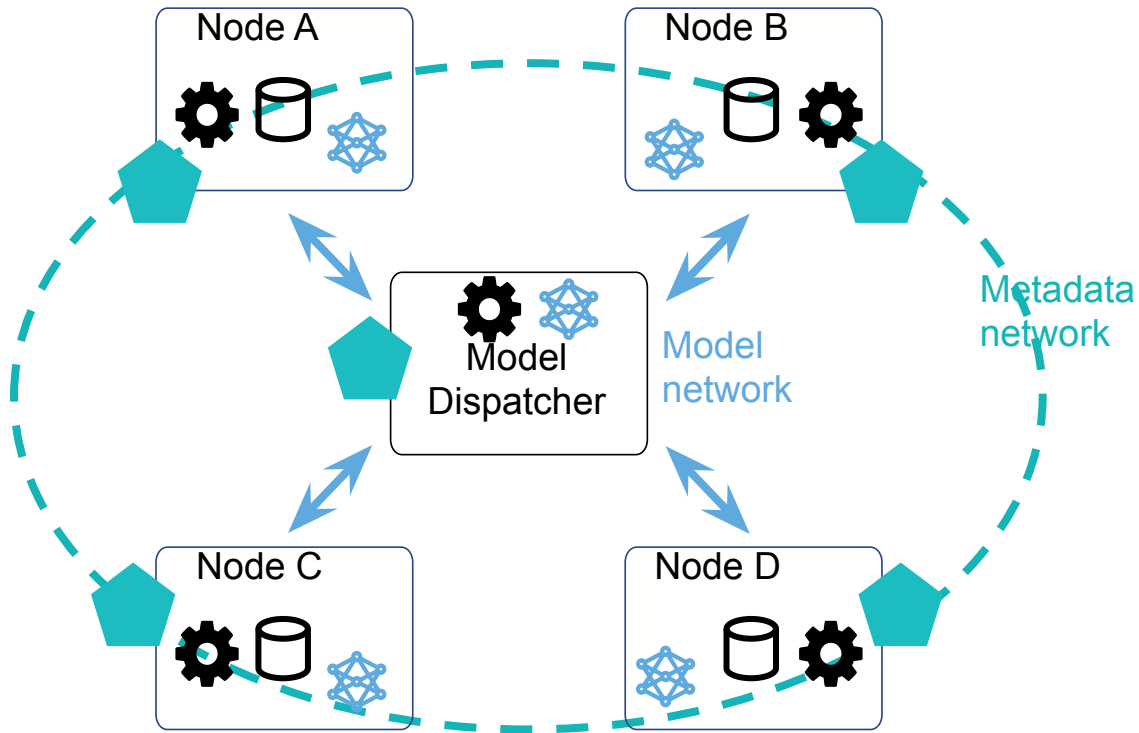
- Algo
    - ML algo and its dependencies

- Models
    - learnt parameters
    - training tasks specification
    - testing tasks specification
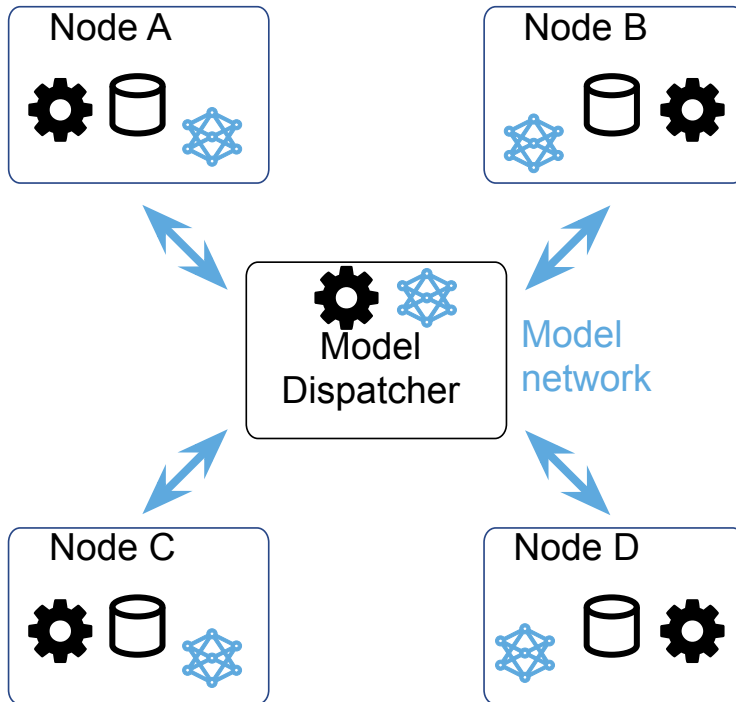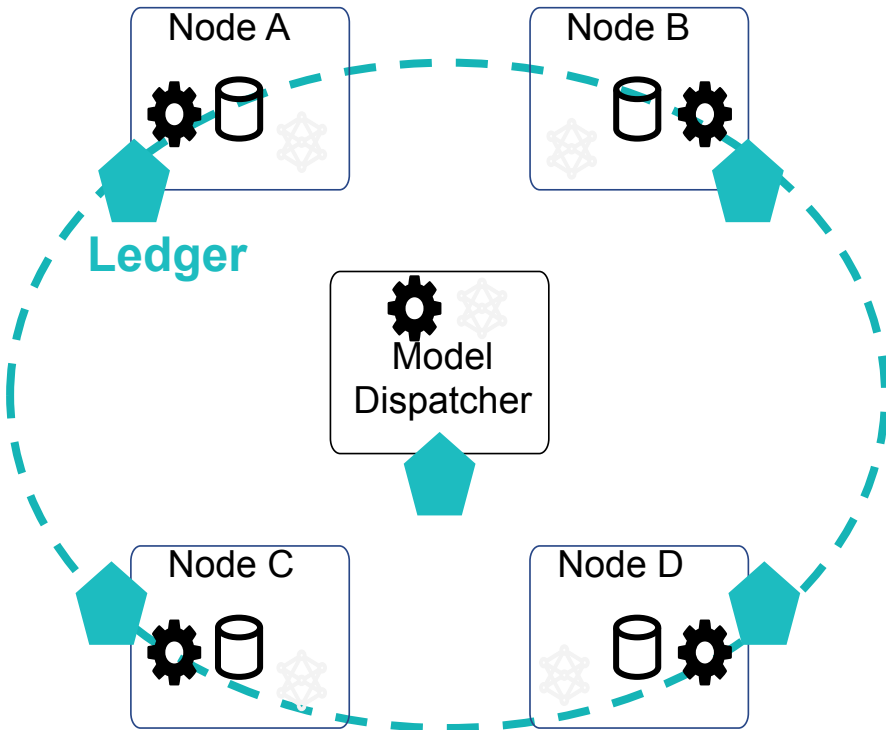    - aggregation tasks specification

OWKIN
# THE TWO NETWORKS



© 2020 OWKIN, INC.

Confidential Information. Do not share without written permission from Owkin, Inc.

12

OWKIN
# THE MODEL NETWORK

Node A

Node B

Model
Dispatcher

Model
network

Node C

Node D

Exchange of model (updates)
through a REST API.

OWKIN

# THE METADATA NETWORK



Node A

Node B

**Ledger**

Model
Dispatcher

Node C

Node D

Distributed Ledger Technology:
Trustless traceability
Learning orchestration
Permissions management

**HYPERLEDGER**
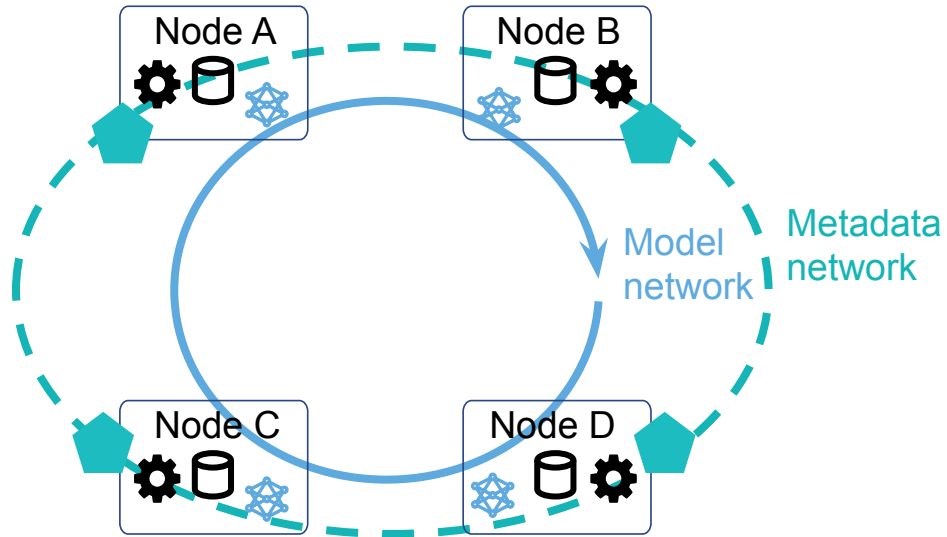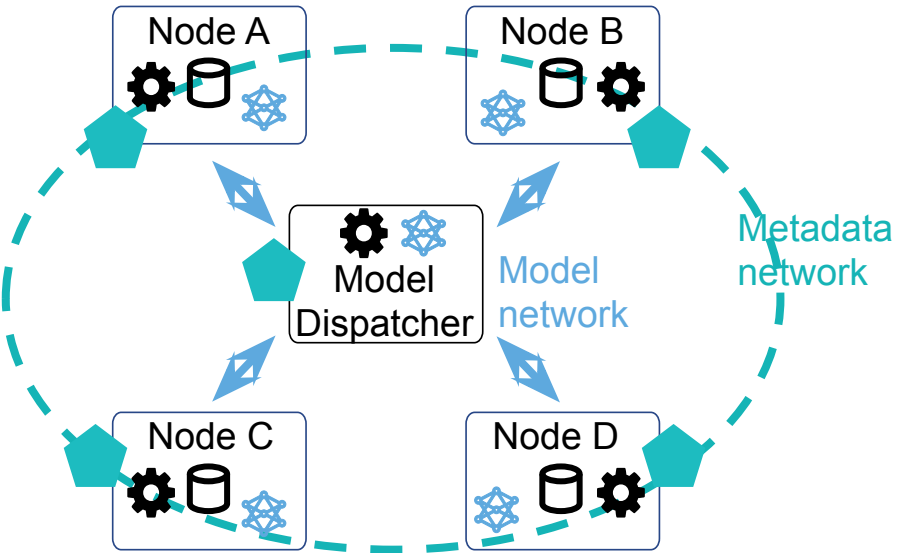FABRIC

OWKIN
# THE METADATA NETWORK



## Metadata of the 4 assets

- Objective
    - Hash of the metrics and url of to request it
    - Id and hash of test data; node that stores them
    - Permissions
- Dataset
    - Hash of the opener and url of to request it
    - Id and hash of data; node that stores them
    - Permissions
- Algo
    - Hash of the algo and url of to request it
    - Node that created it
    - Permissions
- Models
    - Specification of training tasks
    - Specification of aggregation tasks
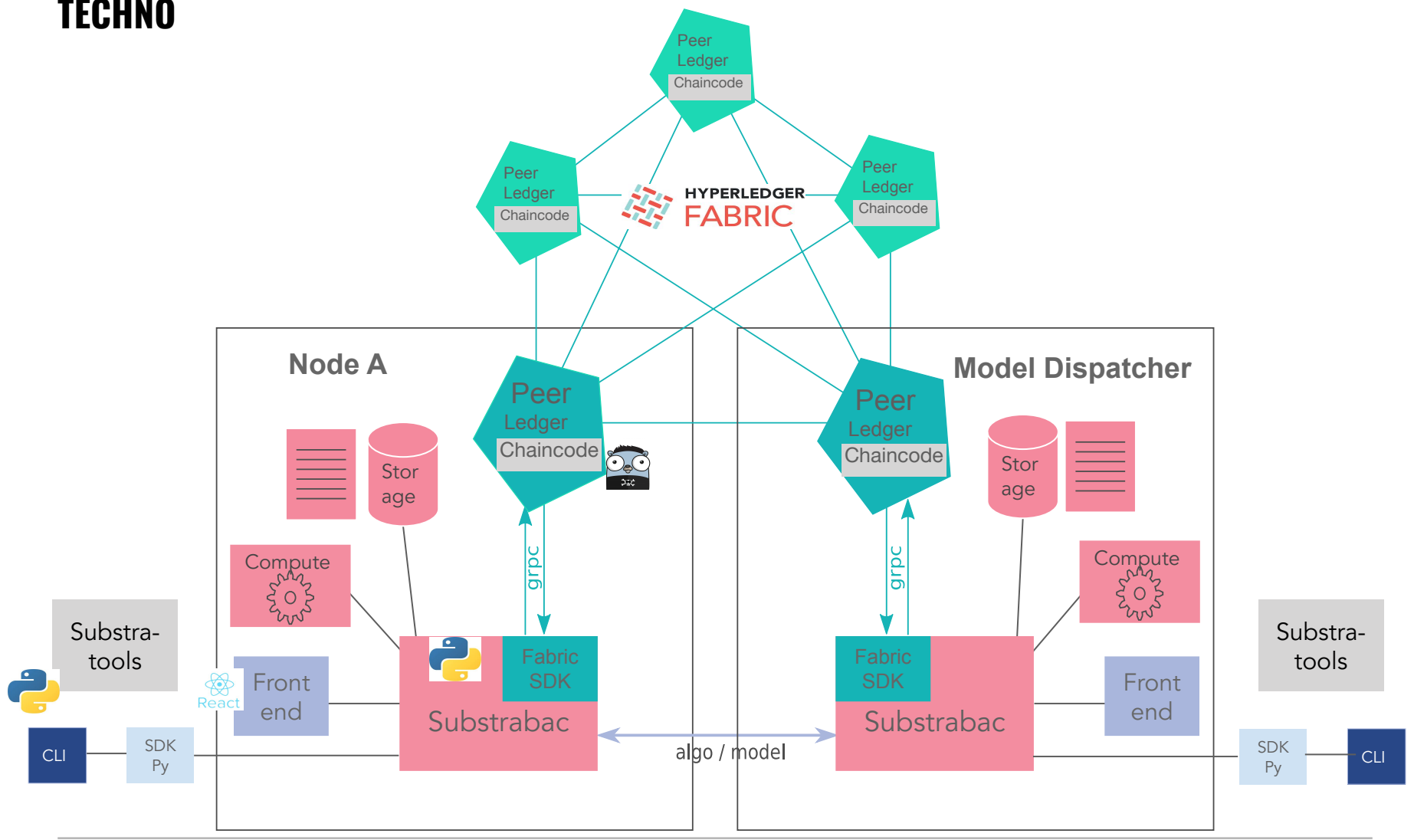    - Specification of evaluation tasks

OWKIN
# NETWORK TOPOLOGIES

OWKIN

# ARCHITECTURE OVERVIEW
## TECHNO

OWKIN
# OWKIN CONNECT INTERFACES

### Frontend
(traceability/perf)



### Command Line Interface

### Python SDK

# OWKIN CONNECT

# COMPUTE PLANS

OWKIN

# FL LEARNING ALGORITHMS



## Different ways to handle a training step and the aggregation of models.

**Algorithm 1** FedSGD

1: Initialize the parameters with some value $\mathbf{w}_0$
2: $C$: proportion of workers chosen at each iteration ($0 < C \leq 1$)
3: $T$: number of aggregation rounds
4: $\{\eta_t\}_{t \geq 1}$: sequence of step sizes
5: **for** $t = 1, \ldots, T$ **do**
6:    Choose uniformly at random $\lfloor CK \rfloor$ distinct edge participants, stored in a set $H = \{k_1, \ldots k_{\lfloor CK \rfloor}\}$.
7:    **for** $k \in H$ **in parallel do**
8:       Edge participant $k$ receives the previous value $\mathbf{w}_{t-1}$ from the central participant
9:       Edge participant $k$ computes the gradient of $\mathcal{L}_k$ with respect to $\mathbf{w}$:

$$\mathbf{g}_{t,k} = \frac{1}{N_k} \sum_{i=1}^{N_k} \nabla_{\mathbf{w}} \mathcal{L}_k(\mathbf{w}_{t-1}) \qquad (3)$$
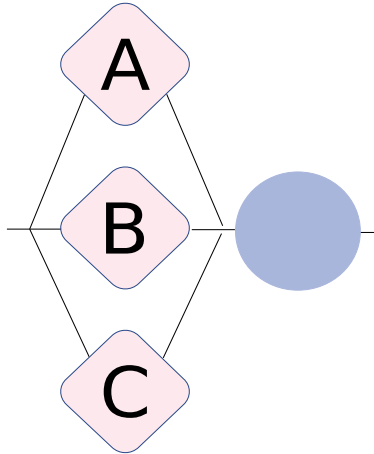
10:       Edge participant $k$ returns local gradient $\mathbf{g}_{t,k}$ to the central server
11:    **end for**
12:    Central computation server aggregates the different gradients:

$$\mathbf{g}_t = \sum_{k \in H} \frac{N_k}{N} \mathbf{g}_{t,k} \qquad (4)$$

13:    Central computation server performs a gradient descent step:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \mathbf{g}_t \qquad (5)$$

14: **end for**
15: **return** $\mathbf{w}_T$

**Algorithm 2** FedAvg

1: Initialize the parameters with some value $\mathbf{w}_0$
2: $C$: proportion of workers chosen at each iteration ($0 < C \leq 1$)
3: $T$: number of aggregation rounds
4: $B$: local batch size
5: $E_l$: number of local epochs
6: $\{\eta_t\}_{t \geq 1}$: sequence of step sizes
7: **for** $t = 1, \ldots, T$ **do**
8:    Choose uniformly at random $\lfloor CK \rfloor$ distinct edge participants, stored in a set $H = \{k_1, \ldots k_{\lfloor CK \rfloor}\}$.
9:    **for** $k \in H$ **in parallel do**
10:       Edge participant $k$ receives the previous value $\mathbf{w}_{t-1}$ from the central participant: $\mathbf{w}_{t,k}^0 \triangleq \mathbf{w}_{t-1}$
11:       **for** $e = 1, \ldots, E_l$ **do**
12:          Initialize value of the parameters for epoch $e$: $\mathbf{w}_{t,k}^e \leftarrow \mathbf{w}_{t,k}^{e-1}$
13:          **for** $b = 1, \ldots, {N_k}/{B}$ **do**
14:             Choose $B$ samples uniformly at random (with or without replacement) from $\mathcal{D}_k$: $s_{i_1}, \ldots, s_{i_B}$
15:             Perform a gradient descent step over the batch

$$\mathbf{w}_{t,k}^e \leftarrow \mathbf{w}_{t,k}^e - \frac{\eta_t}{B} \sum_{j=1}^{B} \nabla_{\mathbf{w}} \ell_{i_j}(\mathbf{w}_{t,k}^e) \qquad (6)$$

16:          **end for**
17:       **end for**
18:       Edge participant $k$ returns the weight updates for the round:

$$\Delta \mathbf{w}_{t,k} = \mathbf{w}_{t,k}^{E_l} - \mathbf{w}_{t-1} \qquad (7)$$

19:    **end for**
20:    Central computation server aggregates the local updates:

$$\Delta \mathbf{w}_t = \sum_{k \in H} \frac{N_k}{N} \Delta \mathbf{w}_{t,k} \qquad (8)$$

21:    Central computation server performs a gradient descent step:

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \Delta \mathbf{w}_t \qquad (9)$$

22: **end for**
23: **return** $\mathbf{w}_T$

McMahan et al, 2017

# Thank you !





## Contact us to learn more:

[https://owkin.com/](https://owkin.com/)

@marini_camille @OWKINscience @Substra_org

Owkin