# TIFeD: a Tiny Integer-based Federated learning algorithm with Direct feedback alignment

Luca Colombo
Politecnico di Milano
Milan, Italy
luca2.colombo@polimi.it

Alessandro Falcetta
Politecnico di Milano
Milan, Italy
alessandro.falcetta@polimi.it

Manuel Roveri
Politecnico di Milano
Milan, Italy
manuel.roveri@polimi.it

## ABSTRACT

Training machine and deep learning models directly on extremely resource-constrained devices is the next challenge in the field of tiny machine learning. The related literature in this field is very limited, since most of the solutions focus only on *on-device* inference or model adaptation through online learning, leaving the training to be carried out on external Cloud services. An interesting technological perspective is to exploit Federated Learning (FL), which allows multiple devices to collaboratively train a shared model in a distributed way. However, the main drawback of state-of-the-art FL algorithms is that they are not suitable for running on tiny devices. For the first time in the literature, in this paper we introduce TIFeD, a Tiny Integer-based Federated learning algorithm with Direct Feedback Alignment (DFA) entirely implemented by using an integer-only arithmetic and being specifically designed to operate on devices with limited resources in terms of memory, computation and energy. Besides the traditional *full-network* operating modality, in which each device of the FL setting trains the entire neural network on its own local data, we propose an innovative *single-layer* TIFeD implementation, which enables each device to train only a portion of the neural network model and opens the door to a new way of distributing the learning procedure across multiple devices. The experimental results show the feasibility and effectiveness of the proposed solution. The proposed TIFeD algorithm, with its *full-network* and *single-layer* implementations, is made available to the scientific community as a public repository.

## CCS CONCEPTS

• **Computing methodologies → Cooperation and coordination**; **Neural networks**; **Artificial intelligence**; **Machine learning**; • **Computer systems organization → Distributed architectures**; **Embedded systems**.

## KEYWORDS

Tiny Machine Learning, Federated Learning, DFA, Deep Learning.

## 1 INTRODUCTION

Tiny Machine Learning (TinyML) is a new emerging research area aiming at bringing Machine Learning (ML) and Deep Learning (DL) algorithms on embedded systems and Internet-of-Things (IoT) devices [1]. The main reason behind this paradigm is that, in recent years, the scientific trend is to move the processing of data as close

as possible to where they are generated. This leads to several advantages. First, it enables decision making directly on devices, hence reducing the latency between data production and data processing and supporting real-time applications. Second, it increases the energy efficiency, given that sending data is more power-consuming than performing computations on device. Third, it enhances privacy and security because sensitive data remain on the IoT units and are not transmitted to remote servers [15].

The main drawback of TinyML, however, is that on-device training of DL models is, in most cases, challenging or even impossible due to the severe technological constraints on memory, computation and energy characterizing IoT units. In the related literature, only a few works, such as [3, 16], have attempted to implement an incremental mechanism, based on transfer learning, to locally adapt the model as new data become available.

To provide on-device training within an IoT system, an interesting technological perspective is to leverage Federated Learning (FL), which is a distributed approach that allows multiple devices to collaboratively train a shared model. Each device performs a local training with its own data, and then sends the model updates to a central server that aggregates the locally computed updates of all the devices in the network [2]. The main drawback of state-of-the-art FL algorithms is that they are designed to operate on edge devices such as smartphones and tablets, which are significantly more powerful from the technological point of view than IoT units. For this reason, traditional FL is not suitable for running on resource-constrained devices, given that: *(i)* the traditional learning procedure with Backpropagation (BP) may not be feasible since it is computationally expensive; *(ii)* the Neural Network (NN) must be shallow and with few parameters; *(iii)* the training dataset has to be small enough to fit into the device's memory.

In this perspective, the aim of this paper is to address the following research question: *how can we train deep and complex machine learning models directly on extremely resource-constrained devices in a federated way?* To the best of our knowledge, we propose for the first time in the literature a Tiny Integer-based Federated learning algorithm with Direct feedback alignment (TIFeD) that introduces the following innovations:

- It enables resource-constrained devices to train a NN without having to rely on external Cloud services by leveraging federated learning. In this way, we can exploit all the advantages of TinyML also during the learning phase and not only for the inference.
- It allows only single layers to be trained, thus reducing the computational and memory demands on the microcontrollers. A direct consequence is that sending weights updates also requires less data since only a portion of the

model needs to be sent to the central server, resulting in energy savings.

- It is entirely implemented through an integer-only arithmetic. In this way, memory consumption can be reduced by using integers less than 32 bits long and, in addition, even devices without floating-point unit (FPU) can train ML and DL models.

The proposed TIFeD algorithm is based on the Direct Feedback Alignment (DFA) [12] learning procedure, which has two main advantages over BP. First, it is much less computationally expensive because each hidden layer is trained independently from the rest of the NN. Second, it allows the implementation of an integer-only arithmetic while avoiding the risk of overflow, which is a major issue when using BP with integer operations.

The paper is organized as follows. Section 2 gives an overview on the related literature. In Section 3, the background on DFA and FL is introduced. Section 4 presents the proposed TIFeD algorithm, while experimental results are shown and discussed in Section 5. Conclusions are finally drawn in Section 6.

## 2 RELATED LITERATURE

This section describes the related literature in the field of federated learning and the available works aiming at implementing the most used FL algorithm, namely Federated Averaging (FedAvg), on tiny devices.

The most popular FL algorithm, introduced by Google in 2016, is Federated Averaging [11]. In FedAvg, each node performs multiple iterations of mini-batch Stochastic Gradient Descent (SGD) with BP to update the local model before sending the gradients to the server. In the central server, received model updates are aggregated together by using an averaging function to combine knowledge learned from different datasets. The main disadvantage of FedAvg is that it has not been designed to operate on resource-constrained devices, as it requires a large amount of memory and high processing power to perform the local learning procedure, and thus it is not suitable in a TinyML scenario. Other works in the field, such as [4, 6, 18], try to optimize the FedAvg algorithm either by improving its efficiency or by using more complex and sophisticated aggregation functions, but without addressing the issue of TinyML-specific design.

On the other hand, [7, 10] are the first papers implementing FedAvg on a device with limited resources, the Arduino Nano 33 BLE Sense board. In particular, Kopparapu et al. [7] exploited transfer learning to solve a binary image classification task, while Llisterri et al. [10] implemented a feed-forward neural network on a simpler keyword spotting task. Given the aforementioned limitations of FedAvg in working on tiny devices, however, both works are only able to train the last fully-connected layer of the considered neural networks.

Differently from the literature, the solution proposed in this paper provides some fundamental aspects for on-device training of ML models, as summarized in Table 1. In particular, TIFeD is specifically designed for operating on resource-constrained devices and it is entirely implemented using an integer-only arithmetic. Moreover, the *single-layer* TIFeD implementation introduces a new

**Table 1: Comparison with existing solutions.**

| Algorithms | Tiny devices-specific design | Integer-only arithmetic | Partial model training |
|---|---|---|---|
| [4] | No | No | No |
| [7] | Yes | No | No |
| [11] | No | No | No |
| **TIFeD** | **Yes** | **Yes** | **Yes** |

operating modality in which each device can train only a portion of the global shared model.

## 3 BACKGROUND

This section introduces the basic concepts on both the DFA training algorithm and the Federated Learning approach.

### 3.1 Direct Feedback Alignment (DFA)

Direct Feedback Alignment [12] is a new learning method for ML and DL models based on the Feedback Alignment (FA) [9] principle. FA showed that the symmetry of weights used in the forward and backward phases, which is required in the BP algorithm, is not necessary. In simpler terms, the network can learn to effectively use fixed and random feedback weights to reduce the error.

The FA training algorithm is based on two insights [9]:

(1) The feedback weights do not need to be exactly equal to the forward weights $W$, but it is sufficient that, for any random matrix $B$ whose elements are uniformly distributed, on average

$$\mathbf{e}^{\top} W B \mathbf{e} > 0 \qquad (1)$$

where $\mathbf{e}$ is the error of the network.

(2) In order to guarantee Eq. (1), instead of adjusting $B$, we can keep it fixed and adjust $W$ accordingly. From a geometrical point of view, Eq. (1) means that the teaching signal used by FA (i.e., $B\mathbf{e}$) lies within 90° of the signal used by BP (i.e., $\mathbf{e}^{\top} W$), that is, $B$ pushes the weights in roughly the same direction as BP. This alignment of signals implies that $B$ acts as $W^{\top}$ and, since $B$ is fixed, the alignment is driven by changes in the forward weights $W$. In this way, even random feedback weights convey useful teaching information throughout the network.

When adopted in deep networks with more than one hidden layer, however, even FA back-propagates the error from the output layer through the upper layers. On the other hand, DFA propagates the output error directly to each of the hidden layers and exploits the FA principle to train them independently of the rest of the neural network.

Let $H$ be the number of hidden layers of a feed-forward NN. Formally, the layer's update directions for FA – denoted by $\delta_{FA}^h$, $h = 1, \ldots, H$ – are calculated as:

$$\delta_{FA}^h = \begin{cases} L' \odot \text{act}'^h(\mathbf{a}^h) & \text{for} \quad h = H \\ \delta_{FA}^{h+1} B^h \odot \text{act}'^h(\mathbf{a}^h) & \text{for} \quad 1 \le h \le H-1 \end{cases}$$

where $L'$ is the derivative of the Loss Function, and $\mathrm{act}'^h(\cdot)$, $\mathbf{a}^h$, and $B^h$ are the derivative of the activation function, the sum of products, and the fixed random weight matrix of layer $h$, respectively.

Concerning DFA, the layer's update directions – denoted by $\delta^h_{DFA}$, $h = 1, \ldots, H$ – are computed as:

$$\delta^h_{DFA} = L' B^h \odot \mathrm{act}'^h(\mathbf{a}^h). \tag{2}$$

Once $\delta^h_{DFA}$, with $h = 1, \ldots, H$, have been computed, the weights updates for each layer are computed as:

$$\begin{aligned} \delta W^h &= -\mathbf{o}^{h-1\,\top} \delta^h_{DFA} \\ \delta \mathbf{b}^h &= -\delta^h_{DFA} \end{aligned} \tag{3}$$

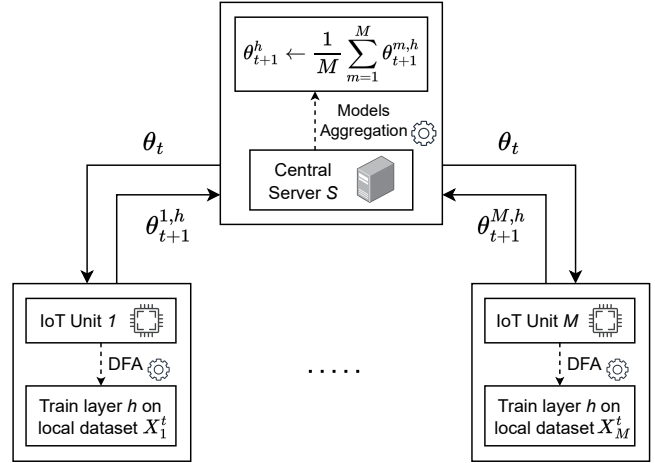where $\mathbf{o}^{h-1}$ is the output of the previous layer in the network.

According to [13], a network trained with DFA converges to the same region of the loss landscape regardless of the initialization of the network weights $W$. This property is guaranteed by the alignment step of the DFA learning procedure, in which the weights $W$ align with the random feedback vectors $B$ used to propagate the network error to each layer. The alignment phase also explains why the random feedback vectors $B$ must be considered as an initialization of the network. Indeed, they drive the learning towards different minima of the loss function, resulting in different performance and accuracies of the final model. DFA has been successfully used in PocketNN [14], a feed-forward neural network tailored for IoT devices.

## 3.2 Federated Learning

Federated learning allows multiple clients to train a shared ML model on decentralized data sources. It addresses the drawbacks of classical ML (i.e. data privacy and security) by distributing the learning procedure across several nodes, thus allowing the data to remain local. Instead of collecting and sending acquired data to an external Cloud service, each device performs a local training by using a shared model as initialization and exchanges only the weights updates with a central server. Local updates are then aggregated by the server to create a global model that will be forwarded to the clients. This process is repeated iteratively until the global model reaches the desired level of accuracy [5, 6].

## 4 THE PROPOSED SOLUTION

This section details the proposed TIFeD algorithm. We emphasize that, as stated in Section 2 and outlined in Table 1, our work introduces the following innovations compared with the literature. First, TIFeD is specifically designed to operate on resource-constrained devices, taking into account their technological limitations in terms of memory, computation and energy. Second, inspired by [14], TIFeD is implemented using an integer-only arithmetic, which enables a reduction in memory consumption by using integers less than 32 bit long and, in addition, allows DL models to be trained on devices without a floating-point unit. Third, the *single-layer* TIFeD implementation introduces a new way of distributing the learning procedure across multiple devices, allowing each IoT unit to train only a portion of the whole NN.



**Figure 1: Overview of the** TIFeD **algorithm operating in a federated learning scenario.**

In particular, Subsection 4.1 gives a general overview on the proposed federated learning scenario. In Subsection 4.2, the transition from DFA to Federated DFA is described, while Subsection 4.3 presents an extension to Federated DFA in which each client node trains only specific layers of the considered neural networks.

## 4.1 Overview

An overview of the proposed TIFeD algorithm operating in a FL scenario is shown in Fig. 1. Without any loss of generality, in the followings a supervised learning image classification task is considered as a reference.

Consider a pervasive system consisting of $M$ IoT units, each of which collects local data $X_m$ through some sensors. The objective is to obtain a ML model trained directly on-device by the IoT units, without relying on external Cloud services. To achieve this goal, the proposed TIFeD algorithm, detailed in the next subsections, exploits FL. The key role in this setup is played by the central server $S$, which is the orchestrator of the entire architecture. At each communication round $t$, the server $S$ sends the global model $\theta_t$ to all the $M$ IoT units, that perform a local computation of the DFA algorithm based on their dataset $X_m^t$. Once the learning procedure is completed, each IoT unit $m$, with $m = 1, \ldots, M$, submits the local updates $\theta_{t+1}^m$ to the server $S$, which is in charge of aggregating the results into a new global model.

In particular, the proposed TIFeD algorithm can operate in two modalities. The first, called *full-network* implementation (detailed in Subsection 4.2), allows each IoT unit to train all the $H$ hidden layers of the DL model. In this case, the result of the computation will be a set containing the weights and bias updates for each layer $h = 1, \ldots, H$. The second, named *single-layer* implementation (explained in Subsection 4.3), enables the IoT units to train only the $h$-th layer of the NN. The result of the computation, in this second scenario, will be a single tuple containing the weights and bias updates related to layer $h$.

---

**Algorithm 1:** *Full-network* TIFeD algorithm implementation. $M$ is the number of devices, $H$ is the number of hidden layers, $bs$ is the mini-batch size, $X_m^t$ is the local dataset at the $t$-th round of client $m$, $e$ is the number of local epochs, and $\eta$ is the learning rate.

---

**central server executes:**
  initialize $B^h$ **foreach** $h = 1, \dots, H$
  **foreach** *round* $t = 1, \dots, T$ **do**
    **foreach** *device* $m = 1, \dots, M$ **in parallel do**
      $\theta_{t+1}^m \longleftarrow$ NodeUpdate$(m, t, \theta_t)$
    **end**
    $\theta_{t+1} \longleftarrow \frac{1}{M} \sum_{m=1}^{M} \theta_{t+1}^m$
  **end**

**NodeUpdate**$(m, t, \theta)$:
  $\mathcal{A} \longleftarrow$ (split $X_m^t$ into mini-batches of size $bs$)
  **foreach** *local epoch* $e$ **do**
    **foreach** *mini-batch* $\alpha \in \mathcal{A}$ **do**
      $\delta\theta \longleftarrow \left\{ [\delta W^1, \delta\mathbf{b}^1], \dots, [\delta W^H, \delta\mathbf{b}^H] \right\}$
      $\theta \longleftarrow \theta - \eta\, \delta\theta$
    **end**
  **end**
  return $\theta_{t+1}^m$ to the central server

---

### 4.2 From DFA to Federated DFA

In the federated learning setting, we consider a fixed set of clients $M$, each with a local dataset $X_m$, and a global model $\theta$ composed of $H$ hidden layers. At the beginning of each round $t$, the central server sends the current model parameter

$$\theta_t = \left\{ [W_t^1, \mathbf{b}_t^1], \dots, [W_t^H, \mathbf{b}_t^H] \right\},$$

being $W_t^h$ and $\mathbf{b}_t^h$ the matrix of weights and the bias vector of layer $h$ respectively, to each of the clients, which perform a local computation of the DFA algorithm. Once the weight updates have been computed by each client $m$ on its dataset $X_m^t$, the new local model parameters $\theta_{t+1}^m$ are sent to the central server, which aggregates all the results and updates its global state with the new model $\theta_{t+1}$.

It is important to note that since we are in an embedded scenario where devices have a limited amount of memory (e.g., the Arduino Nano 33 BLE Sense board has 256kB of SRAM), the local dataset $X_m$ cannot be arbitrarily big, because it must fit into the device's memory. This is a major limitation compared to traditional FL, where learning procedures benefit from very large and varied datasets.

More in detail, the *full-network* TIFeD algorithm implementation, summarized in Algorithm 1, can be divided into two blocks: the central server component and the IoT units component. Concerning the former part, executed by the **central server**, it works as follows:

(1) The node initializes, for each hidden layer $h$ of the network, the matrix of feedback weights $B^h$ by sampling values from a uniform distribution, and sends them to each client $m$.

---

**Algorithm 2:** *Single-layer* TIFeD algorithm implementation. $M$ is the number of devices, $H$ is the number of hidden layers, $bs$ is the mini-batch size, $X_m^t$ is the local dataset at the $t$-th round of client $m$, $e$ is the number of local epochs, and $\eta$ is the learning rate.

---

**central server executes:**
  initialize $B^h$ **foreach** $h = 1, \dots, H$
  $s \longleftarrow \lfloor \frac{M}{H} \rfloor$
  **foreach** *round* $t = 1, \dots, T$ **do**
    $\mathcal{S}_h \longleftarrow$ (split the clients into $H$ groups of size $s$)
    **foreach** *group* $h = 1, \dots, H$ **in parallel do**
      **foreach** *device* $m \in \mathcal{S}_h$ **in parallel do**
        $\theta_{t+1}^{m,h} \longleftarrow$ NodeUpdate$(m, t, h, \theta_t)$
      **end**
      $\theta_{t+1}^h \longleftarrow \frac{1}{s} \sum_{m=1}^{s} \theta_{t+1}^{m,h}$
    **end**
    $\theta_{t+1} \longleftarrow \left\{ \theta_{t+1}^1, \dots, \theta_{t+1}^H \right\}$
  **end**

**NodeUpdate**$(m, t, h, \theta)$:
  $\mathcal{A} \longleftarrow$ (split $X_m^t$ into mini-batches of size $bs$)
  **foreach** *local epoch* $e$ **do**
    **foreach** *mini-batch* $\alpha \in \mathcal{A}$ **do**
      $\delta\theta^h \longleftarrow \left\{ [\delta W^h, \delta\mathbf{b}^h] \right\}$
      $\theta^h \longleftarrow \theta^h - \eta\, \delta\theta^h$
    **end**
  **end**
  return $\theta_{t+1}^{m,h}$ to the central server

---

(2) For each round $t = 1, \dots, T$, the node sends the current model parameters $\theta_t$ to the clients in order to let them perform local updates using DFA.

(3) After receiving the $M$ local updates $\theta_{t+1}^m$ from the clients, the node upgrades the global model $\theta_{t+1}$ by aggregating the results with the average function:

$$\theta_{t+1} \longleftarrow \frac{1}{M} \sum_{m=1}^{M} \theta_{t+1}^m.$$

Regarding the **NodeUpdate** function, executed by each IoT unit, it works as follows:

(1) The node splits its local part of the dataset $X_m^t$ (relative to the current round $t$) into mini-batches $\alpha$ of size $bs$.

(2) For each local epoch $e$ and for each mini-batch $\alpha$, the node computes the weights updates through the DFA learning algorithm, as shown in Eq. (3), and applies them to the current model parameters $\theta_t$.

(3) Once the computation is completed, the node sends the new model parameters $\theta_{t+1}$ back to the central server.

The models aggregation can be done at different points: after each mini-batch $\alpha$, after the entire buffer $X_m^t$, or after the number of training passes $e$ each client makes on the current buffer $X_m^t$. The choice of the aggregation point represents a trade-off between

the accuracy of the final model and the number of required communication rounds, which results into higher power consumption. A detailed analysis about this point is provided in the experimental results in Section 5.

### 4.3 Single layers training

One of the main advantages of the DFA learning algorithm is the ability to train each NN layer independently. Indeed, as shown in Eq. (2), the layer's update direction formula $\delta_{DFA}^h$ is not recursive and, for each hidden layer $h$, depends only on the derivative of the Loss Function $L'$, the fixed matrix of feedback weights $B^h$ and the derivative of the activation function $act'^h(\cdot)$. This means that, unlike BP, the network error is propagated directly to each hidden layer and does not have to traverse the entire network.

By exploiting this feature, it is possible to design a federated algorithm in which each IoT unit, instead of training an entire network on its local dataset, trains only a single layer of the global model. This results in two benefits: *(i)* the algorithm is more lightweight in terms of computational and memory demands and can therefore run on devices with tighter resource constraints, and *(ii)* the communication between the devices and the central server requires the exchange of a smaller amount of data since only a portion of the network (i.e., only the updates $\theta_{t+1}^{m,h}$ of the $h$-th layer) needs to be sent, resulting in energy savings. On the other hand, the drawbacks are that some overhead is added to the central server that has to decide which device trains which layer, and that there is a loss in the accuracy of the final global model. More accurate analyses are provided in Section 5.

The *single-layer* TIFeD algorithm implementation is shown in Algorithm 2 and works as follows:

- The central server also defines, for each round $t$, the set $S_h$ which contains a random splitting of the clients into $H$ different groups of size $s \longleftarrow \left\lfloor \frac{M}{H} \right\rfloor$. Each group is in charge of training only the $h$-th layer of the network.
- The IoT units compute the weights update related only to the $h$-th layer through the DFA learning algorithm, apply it to the current model parameters $\theta_t^h$, and send the new model parameters $\theta_{t+1}^h$ back to the central server.
- After receiving the $M$ local updates $\theta_{t+1}^{m,h}$ from the clients, the server updates the global model $\theta_{t+1}$ by aggregating the results with the average function, layer-by-layer:

$$\theta_{t+1} \longleftarrow \left\{ \theta_{t+1}^1, \ldots, \theta_{t+1}^H \right\},$$

$$\theta_{t+1}^h \longleftarrow \frac{1}{s} \sum_{m=1}^{s} \theta_{t+1}^{m,h}$$

for $h = 1, \ldots, H$.

## 5 EXPERIMENTAL RESULTS

In this section, the proposed TIFeD algorithm with its *full-network* and *single-layer* implementations are evaluated from different points of view. The code of the following experiments, written in Python, can be found in the code repository[1].

---

[1]https://github.com/AI-Tech-Research-Lab/TIFeD

**Table 2: NNs architecture considered in this experimental section. The rescaled Piecewise Linear Approximation (PLA) of the *tanh* function [14] was used after each FC layer, whereas after the Conv2D layers the *ReLU* activation function was employed. The kernel size for the Conv2D layers is $(3 \times 3)$, while for the MaxPooling layer is $(2 \times 2)$. * Only the classifier is trained using the proposed TIFeD algorithm.**

| Model | Layers |
|---|---|
| NN-1 | $FC(28 * 28) \rightarrow FC(200) \rightarrow FC(10)$ |
| CNN-1 * | $Conv2D(4) \rightarrow MaxPool \rightarrow Conv2D(8) \rightarrow$ $MaxPool \rightarrow FC(200) \rightarrow FC(50) \rightarrow FC(10)$ |
| CNN-2 * | $Conv2D(128) \rightarrow MaxPool \rightarrow Conv2D(256) \rightarrow$ $MaxPool \rightarrow Conv2D(512) \rightarrow Conv2D(256) \rightarrow$ $MaxPool \rightarrow FC(256) \rightarrow FC(128) \rightarrow FC(10)$ |

In the following experiments, we consider a fixed set of clients $M$, each with a local dataset $X_m$ uniformly sampled from the training set $X$. It is important to note that since we are in a tiny scenario where devices have a limited amount of memory (e.g., the Arduino Nano 33 BLE Sense board has 256kB of SRAM), the local dataset $X_m$ cannot be processed all at once as if it were entirely available to the device, because it could not fit into the memory of the device. Therefore, the local dataset $X_m$ is further divided into smaller portions of length buff_len and, at each round $t$, only one portion of the dataset of size buff_len is considered for training (i.e., the $t$-th buffer $X_m^t$).
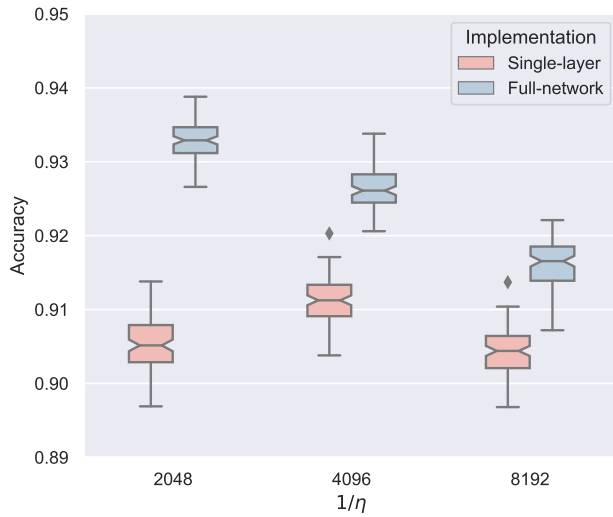
### 5.1 Datasets

A brief description of the datasets employed in our study is given in the following.

*5.1.1 MNIST.* The first considered dataset is the MNIST (Modified National Institute of Standards and Technology) database [19], a collection of handwritten digits commonly used for training and testing ML models able to recognize images. The elements are grey-scale images of size $28 \times 28$. The dataset is composed of 60000 training images and 10000 testing images from 10 different classes, each one representing a digit from 0 to 9.
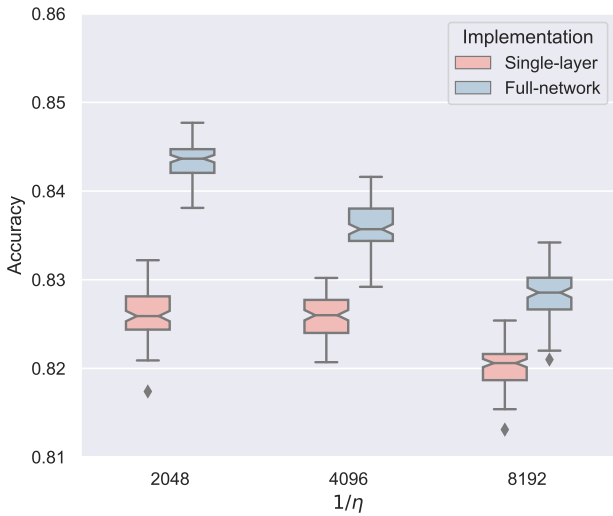
*5.1.2 FashionMNIST.* The FashionMNIST database [17] was created in 2017 as an improvement of the original MNIST dataset for benchmarking ML algorithms. FashionMNIST consists in gray-scale images of dimension $28 \times 28$ of Zalando's articles. The training dataset contains 60000 elements, while the testing dataset includes 10000 elements. Each image is then associated with a label from 10 different classes.

*5.1.3 CIFAR10.* Lastly, the CIFAR10 (Canadian Institute for Advanced Research) dataset [8] is considered. CIFAR10 is a collection of colour images representing objects with size $32 \times 32$ belonging to 10 different classes. The dataset is divided into 50000 training samples and 10000 test samples.
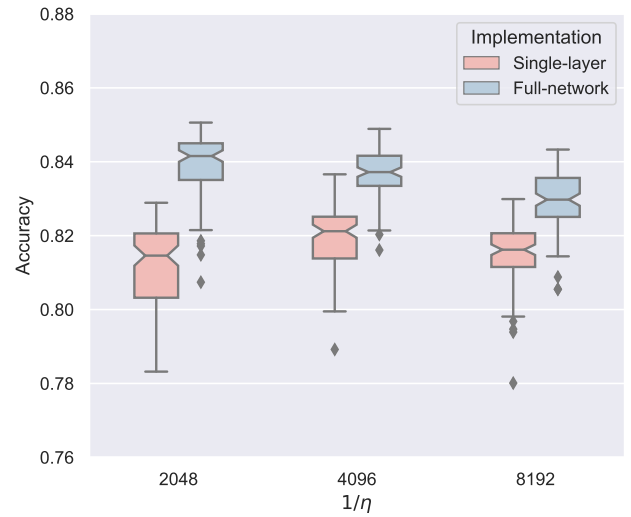
**(a) MNIST dataset**



**(b) FashionMNIST dataset**

**Figure 2: Test accuracies of NN-1 trained with the proposed *full-network* and *single-layer* TIFeD implementations in a federated setting composed of $M = 8$ worker nodes. NN-1 hyperparameters: length of the local dataset $X_m = 7.5k$ images, buff_len $= 50$, $bs = 25$, and $e = 5$. Three different values of the learning rate $\frac{1}{\eta}$ were considered. Results are computed for 100 different initialization of the DFA random-feedback weights $B^h$.**

## 5.2 Feed-forward neural network training

The goal of the first experiment is to show that the proposed TIFeD algorithm is able to train a feed-forward neural network from scratch. To this end, we trained a neural network composed of two fully-connected layers (NN-1 in Table 2) in a federated setting with $M = 8$ worker nodes on the MNIST and FashionMNIST



**Figure 3: Test accuracies of CNN-1 trained, by exploiting transfer learning, with the proposed *full-network* and *single-layer* TIFeD implementations in a federated setting composed of $M = 8$ worker nodes on the FashionMNIST dataset. CNN-1 hyperparameters: length of the local dataset $X_m = 7.5k$ images, buff_len $= 50$, $bs = 25$, and $e = 5$. Three different values of the learning rate $\frac{1}{\eta}$ were considered. Results are computed for 100 different initialization of the DFA random-feedback weights $B^h$.**

datasets with the following set of hyperparameters: length of the local dataset $X_m = 7.5k$ images, buff_len $= 50$, $bs = 25$, and $e = 5$. We evaluated the test accuracy for both the *full-network* and *single-layer* implementations with three different values of the learning rate $\frac{1}{\eta}$ and for 100 different initialization of the DFA random-feedback weights $B^h$.

As we can see in Figure 2, the $M$ worker nodes are able to train a shared model initialized with zero-valued weights by using the proposed TIFeD algorithm. It is important to note that the *single-layer* implementation performs worse than the *full-network* implementation. This can be easily explained by the fact that, in the former case, each FC layer of the NN was trained by 4 out of 8 nodes, while, in the latter case, each NN layer is trained by all the 8 nodes in the federated setting. On the other hand, the computational, memory, and energy demands of the nodes are halved in the *single-layer* implementation compared to the *full-network* implementation, given that weight updates have to be computed and sent for only one layer $h$. As a consequence, the *single-layer* implementation is more suitable when devices are extremely constrained in terms of resources.

## 5.3 Transfer learning for convolutional neural networks training

In a TinyML environment, the ability to adapt a general and complex ML model to the specific application scenario is a fundamental aspect. This second experiment aims to prove that, in addition to feed-forward NNs, it is possible to use TIFeD to train the classifier
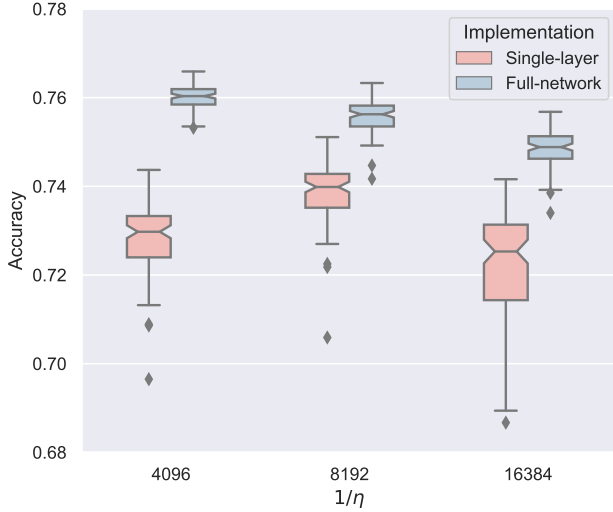
**Figure 4: Test accuracies of CNN-2 trained, by exploiting transfer learning, with the proposed *full-network* and *single-layer* TIFeD implementations in a federated setting composed of $M = 8$ worker nodes on the CIFAR10 dataset. CNN-2 hyperparameters: length of the local dataset $X_m = 6.25k$ images, buff_len = 50, $bs = 25$, and $e = 10$. Three different values of the learning rate $\frac{1}{\eta}$ were considered. Results are computed for 100 different initialization of the DFA random-feedback weights $B^h$.**
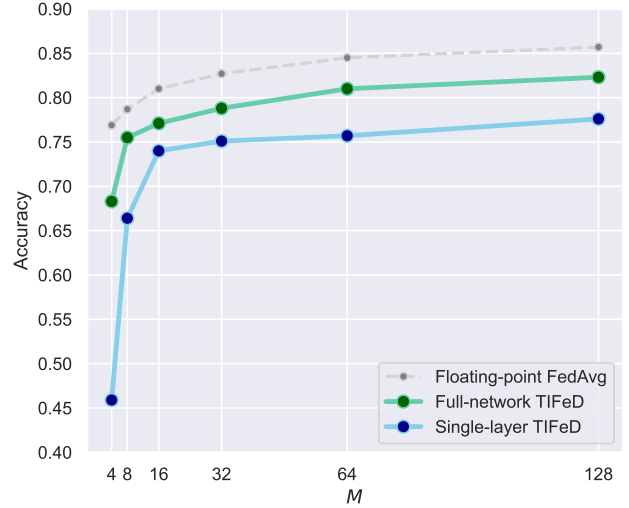


**Figure 5: Analysis of the behaviour of *full-network* TIFeD and *single-layer* TIFeD as the number of worker nodes $M$ increases. Floating-point FedAvg is taken as a reference, but we emphasize that the goal is not to compare their final accuracies, since the two algorithms are tailored for two different application scenarios. Training of CNN-1, by exploiting transfer learning, on the FashionMNIST dataset. FedAvg hyperparameters: $X_m = 100$ images, buff_len = 10, $bs = 10$, $e = 10$, and $\eta = 0.1$. TIFeD hyperparameters: $X_m = 100$ images, buff_len = 20, $bs = 10$, $e = 10$, and $\frac{1}{\eta} = 2048$. Results are computed for 100 different initialization of the DFA random-feedback weights $B^h$ for TIFeD and 100 different network initialization for FedAvg.**

of a Convolutional Neural Network (CNN) that exploits transfer learning.

First, we reconsidered the experiment presented in Section 5.2 on the FashionMNIST dataset. We trained CNN-1 and quantized the feature extractor weights, which are then used by each worker node $m$. At this point, by exploiting TL, the FC layers of CNN-1 were re-trained from scratch directly on the worker nodes, using the same set of hyperparameters as in the experiment of Section 5.2. Similarly to the previous experiment, as can be seen from the results in Figure 3, both the implementations of TIFeD are able to train the considered model.

Then, by adopting the same procedure, we trained CNN-2 on the more complex CIFAR10 dataset. The following set of hyperparameters was used: length of the local dataset $X_m = 6.25k$ images, buff_len = 50, $bs = 25$, and $e = 10$. The results, provided in Figure 4, show that TIFeD is capable of making CNN-2 learn this more challenging image classification task, that a simpler feed-forward NN would not be able to solve.

## 5.4 Exploring TIFeD as the number of worker nodes increases

In this third experiment, the behaviour of TIFeD as the number of worker nodes $M$ increases is analyzed. As a reference, we consider the state-of-the-art FL algorithm, i.e., Federated Averaging. We emphasize that this experiment is not intended to compare the final accuracy achieved by our TIFeD algorithm with respect to FedAvg.

In fact, as previously mentioned, they are tailored for two deeply different application scenarios: TIFeD is specifically designed to run on resource-constrained devices, while FedAvg would not be able to be executed on tiny devices.

We trained the classifier of CNN-1 on the FashionMNIST dataset for different values of $M$. In particular, we set a fixed number of 100 images per worker and trained both FedAvg, which uses 32-bit long floating-point values, and the *full-network* and *single-layer* implementations of the TIFeD algorithm, which operate with 16-bit long integers and an integer-only arithmetic. The experiment was performed using the following set of hyperparameters for both algorithms: length of the local dataset $X_m = 100$ images, $bs = 10$, and $e = 10$. To take into account the lower memory consumption of TIFeD with respect to FedAvg, we set buff_len = 20 for the former, buff_len = 10 for the latter.

The results in Figure 5 show that the *full-network* and *single-layer* TIFeD implementations behave exactly as FedAvg when increasing the number of worker nodes $M$. It should be noted that for the same number of nodes $M$, FedAvg performs twice as many communication rounds than the *full-network* TIFeD implementation, since its buffer length is half as large, resulting in a longer and more energy-consuming training.

**Table 3: Impact of the aggregation point within the *full-network* and *single-layer* TIFeD implementations in terms of test accuracy and number of communication rounds $t$. Training of CNN-1, by exploiting transfer learning, in a federated setting composed of $M = 128$ worker nodes on the Fashion-MNIST dataset. CNN-1 hyperparameters: length of the local dataset $X_m = 100$ images, buff_len $= 20$, $bs = 10$, $e = 10$, and $\frac{1}{\eta} = 2048$. Results are computed for 100 different initialization of the DFA random-feedback weights $B^h$.**

| TIFeD Implementation | Aggregation Point | Test Accuracy | Communication Rounds |
|---|---|---|---|
| *Full-network* | Mini-batch $\alpha$ | 0.819 | 100 |
| | Buffer $X_m^t$ | 0.841 | 50 |
| | Epochs $e$ | 0.823 | 5 |
| *Single-layer* | Mini-batch $\alpha$ | 0.780 | 100 |
| | Buffer $X_m^t$ | 0.791 | 50 |
| | Epochs $e$ | 0.776 | 5 |

## 5.5 The impact of the aggregation point within the TIFeD algorithm

Lastly, an experiment to evaluate the impact of the aggregation point within the TIFeD algorithm, is proposed. As explained in Section 4.2, the aggregation of the weights updates by the central server can occur at three different points: after each mini-batch $\alpha$, after the entire buffer $X_m^t$, or after the number of training epochs $e$. To this end, we performed the learning of the CNN-1 classifier with both the implementations of TIFeD on the FashionMNIST dataset, using the following set of hyperparameters: $M = 128$, $X_m = 100$ images, buff_len $= 20$, $bs = 10$, $e = 10$, and $\frac{1}{\eta} = 2048$.

As summarized in Table 3, the choice of the aggregation point represents a trade-off between the final test accuracy and the number of required communication rounds for both the *full-network* and *single-layer* implementations. In particular, the *Epochs* aggregation point represents the one with the lowest number of required rounds, but it results in a final accuracy drop of 1.5% and 1.8% compared to that provided by the *Buffer* aggregation point. On the other hand, the latter requires $e$ times more communication rounds than the former, resulting in more energy consumption by the tiny devices.

## 6 CONCLUSIONS

The aim of this paper was to present a novel federated learning algorithm specifically designed for training ML and DL models directly on extremely resource-constrained devices. Specifically, we introduced TIFeD, a Tiny Integer-based Federated learning algorithm with Direct feedback alignment, with its two variants: the *full-network* implementation, which allows each node in the federated setting to train all the fully-connected layers of the neural network, and the *single-layer* implementation, that enables the devices to train only a single layer of the entire NN, resulting in a more lightweight algorithm in terms of computation, memory and energy required. The experimental results show the feasibility

and effectiveness of the proposed solution. Future works will consider new models as well as further optimizations in the learning procedure.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Colby R Banbury, Vijay Janapa Reddi, Max Lam, William Fu, Amin Fazel, Jeremy Holleman, Xinyuan Huang, Robert Hurtado, David Kanter, Anton Lokhmotov, et al. 2020. Benchmarking tinyml systems: Challenges and direction. *arXiv preprint arXiv:2003.04821* (2020).
[2] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloe Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMahan, et al. 2019. Towards federated learning at scale: System design. *Proceedings of machine learning and systems* 1 (2019), 374–388.
[3] Simone Disabato and Manuel Roveri. 2020. Incremental on-device tiny machine learning. In *Proceedings of the 2nd International workshop on challenges in artificial intelligence and machine learning for internet of things*. 7–13.
[4] Jack Goetz, Kshitiz Malik, Duc Bui, Seungwhan Moon, Honglei Liu, and Anuj Kumar. 2019. Active federated learning. *arXiv preprint arXiv:1909.12641* (2019).
[5] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. 2021. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning* 14, 1–2 (2021), 1–210.
[6] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. 2016. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492* (2016).
[7] Kavya Kopparapu and Eric Lin. 2021. TinyFedTL: Federated transfer learning on tiny devices. *arXiv preprint arXiv:2110.01107* (2021).
[8] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. 2009. Cifar-10 (canadian institute for advanced research). 2009. *URL http://www. cs. toronto. edu/kriz/cifar. html* 5 (2009).
[9] Timothy P Lillicrap, Daniel Cownden, Douglas B Tweed, and Colin J Akerman. 2014. Random feedback weights support learning in deep neural networks. *arXiv preprint arXiv:1411.0247* (2014).
[10] Nil Llisterri Giménez, Marc Monfort Grau, Roger Pueyo Centelles, and Felix Freitag. 2022. On-device training of machine learning models on microcontrollers with federated learning. *Electronics* 11, 4 (2022), 573.
[11] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*. PMLR, 1273–1282.
[12] Arild Nøkland. 2016. Direct feedback alignment provides learning in deep neural networks. *Advances in neural information processing systems* 29 (2016).
[13] Maria Refinetti, Stéphane d'Ascoli, Ruben Ohana, and Sebastian Goldt. 2021. Align, then memorise: the dynamics of learning with feedback alignment. In *International Conference on Machine Learning*. PMLR, 8925–8935.
[14] Jaewoo Song and Fangzhen Lin. 2022. PocketNN: Integer-only Training and Inference of Neural Networks via Direct Feedback Alignment and Pocket Activations in Pure C++. *arXiv preprint arXiv:2201.02863* (2022).
[15] John A Stankovic. 1996. Real-time and embedded systems. *ACM Computing Surveys (CSUR)* 28, 1 (1996), 205–208.
[16] Prahalathan Sundaramoorthy, Gautham Krishna Gudur, Manav Rajiv Moorthy, R Nidhi Bhandari, and Vineeth Vijayaraghavan. 2018. Harnet: Towards on-device incremental learning using deep ensembles on constrained devices. In *Proceedings of the 2nd International Workshop on Embedded and Mobile Deep Learning*. 31–36.
[17] Han Xiao, Kashif Rasul, and Roland Vollgraf. 2017. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747* (2017).
[18] Kai Yang, Tao Jiang, Yuanming Shi, and Zhi Ding. 2020. Federated learning via over-the-air computation. *IEEE Transactions on Wireless Communications* 19, 3 (2020), 2022–2035.
[19] LeCun Yann. 1998. The mnist database of handwritten digits. *R* (1998).