# Model-Agnostic Federated Learning

Gianluca Mittone[1][0000−0002−1887−6911], Walter Riviera[2][0000−0001−5292−7594], Iacopo Colonnelli[1][0000−0001−9290−2017], Robert Birke[1][0000−0003−1144−3707], and Marco Aldinucci[1][0000−0001−8788−0829]

[1] University of Turin, Turin, Italy
{gianluca.mittone, iacopo.colonnelli, robert.birke, marco.aldinucci}@unito.it
[2] University of Verona, Verona, Italy
walter.riviera@univr.it

**Abstract.** Since its debut in 2016, Federated Learning (FL) has been tied to the inner workings of Deep Neural Networks (DNNs); this allowed its development as DNNs proliferated but neglected those scenarios in which using DNNs is not possible or advantageous. The fact that most current FL frameworks only support DNNs reinforces this problem. To address the lack of non-DNN-based FL solutions, we propose MAFL (Model-Agnostic Federated Learning). MAFL merges a model-agnostic FL algorithm, AdaBoost.F, with an open industry-grade FL framework: Intel® OpenFL. MAFL is the first FL system not tied to any machine learning model, allowing exploration of FL beyond DNNs. We test MAFL from multiple points of view, assessing its correctness, flexibility, and scaling properties up to 64 nodes of an HPC cluster. We also show how we optimised OpenFL achieving a 5.5x speedup over a standard FL scenario. MAFL is compatible with x86-64, ARM-v8, Power and RISC-V.

**Keywords:** Machine Learning · Federated Learning · Federated AdaBoost · Software Engineering

## 1 Introduction

Federated Learning (FL) is a Machine Learning (ML) technique that has gained tremendous popularity in the last years [9]: a shared ML model is trained without ever exchanging the data owned by each party or requiring it to be gathered in one common computational infrastructure. The popularity of FL caused the development of a plethora of FL frameworks, e.g., Flower [4], FedML [7], and HPE Swarm Learning [23] to cite a few. These frameworks only support one ML model type: Deep Neural Networks (DNNs). While DNNs have shown unprecedented results across a wide range of applications, from image recognition [11] to natural language processing [22], from drug discovery [24] to fraud detection [10], they are not the best model for every use case. DNNs require massive amounts of data, which collecting and eventually labelling is often prohibitive; furthermore, DNNs are not well-suited for all types of data. For example, traditional ML models can offer a better performance-to-complexity ratio on tabular data

than DNNs [17]. DNNs also behave as black-box, making them undesirable when the model's output has to be explained [8]. Lastly, DNNs require high computational resources, and modern security-preserving approaches, e.g. [16, 21], only exacerbate this issues [14].

We propose the open-source **MAFL**[1] (*Model-Agnostic Federated Learning*) framework to alleviate these problems. MAFL leverages *Ensemble Learning* to support and aggregate ML models independently from their type. Ensemble Learning exploits the combination of multiple *weak learners* to obtain a single *strong learner*. A weak learner is a learning algorithm that only guarantees performance better than a random guessing model; in contrast, a strong learner provides a very high learning performance (at least on the training set). Since weak learners are not bound to be a specific ML model, Ensemble Learning techniques can be considered *model-agnostic*. We adopt the AdaBoost.F algorithm [18], which leverages the AdaBoost algorithm [6] and adapts it to the FL setting, and we marry it with an open-source industry-grade FL platform, i.e., Intel$^{®}$ OpenFL [5]. To our knowledge, MAFL is the first and only model-agnostic FL framework available to researchers and industry at publication.

The rest of the paper introduces the basic concepts behind MAFL. We provide implementation details underlying its development, highlight the challenges we overcame, and empirically assess our approach from the computational performances and learning metrics points of view. To summarise, the contributions of this paper are the following:

- we introduce MAFL, the first FL software able to work with any supervised ML model, from heavy DNNs to lightweight trees;
- we describe the architectural challenges posed by a model-agnostic FL framework in detail;
- we describe how Intel$^{®}$ OpenFL can be improved to boost computational performances;
- we provide an extensive empirical evaluation of MAFL to showcase its correctness, flexibility, and performance.

## 2   Related Works

*FL* [15] usually refers to a centralised structure in which two types of entities, a single *aggregator* and multiple *collaborators*, work together to solve a common ML problem. A FL framework orchestrate the federation by distributing initial models, collecting the model updates, merging them according to an aggregation strategy, and broadcasting back the updated model. FL requires a *higher-level software infrastructure* than traditional ML flows due to the necessity of exchanging model parameters quickly and securely. Model training is typically delegated to de-facto standard (deep) ML frameworks, e.g., PyTorch and TensorFlow.

Different *FL frameworks* are emerging. Riviera [19] provides a compelling list of 36 open-source tools ranked by community adoption, popularity growth, and

---

[1] https://github.com/alpha-unito/Model-Agnostic-FL

feature maturity, and Beltrán [3] reviews 16 FL frameworks, identifying only six as mature. All of the surveyed frameworks support supervised training of DNNs, but only FATE [12], IBM-Federated [13], and NVIDIA FLARE [20] offer support for a few different ML models, mainly implementing federated K-means or Extreme Gradient Boosting (XGBoost): this is due to the problem of defining a model-agnostic aggregation strategy. DNNs' client updates consist of tensors (mainly weights or gradients) that can be easily serialised and mathematically combined (e.g., averaged), as are also the updates provided by federated K-means and XGBoost. This assumption does not hold in a model-agnostic scenario, where the serialisation infrastructure and the aggregation mechanism have to be powerful enough to accommodate different update types. A truly model-agnostic aggregation strategy should be able to aggregate not only tensors, but also complex objects like entire ML model. AdaBoost.F is capable of doing that. Section 3 delves deeper into the state-of-the-art of federated ensemble algorithms.

As a base for developing MAFL, we chose a mature, open-source framework supporting only DNNs: Intel® OpenFL [5]. The reason for this choice is twofold: (i) its structure and community support; and (ii) the possibility of leveraging the existing ecosystem by maintaining the same use and feel. Section 4 delves into the differences between plain OpenFL and its MAFL extension, showing how much DNN-centric a representative modern FL framework can be.

## 3  Model-agnostic Federated Algorithms

None of the frameworks mentioned in Sec. 2 supports model-agnostic FL algorithms, i.e., they cannot handle different ML models seamlessly. The reason is twofold. On the one hand, modern FL frameworks still try to achieve sufficient technical maturity, rather than adding new functionalities. On the other hand, model-agnostic federated algorithms are still new and little investigated.

Recently, [18] proposed three federated versions of AdaBoost: *DistBoost.F*, *PreWeak.F*, and *AdaBoost.F*. All three algorithms are model-agnostic due to their inherent roots in AdaBoost. Following the terminology commonly used in ensemble learning literature, we call *weak hypothesis* a model learned at each federated round and *strong hypothesis* the final global model produced by the algorithms. The general steps of an AdaBoost-based FL algorithm are the following:

1. The aggregator receives the dataset size $N$ from each collaborator and sends them an initial version of the weak hypothesis.
2. The aggregator receives the weak hypothesis $h_i$ from each collaborator and broadcasts the entire hypothesis space to every collaborator.
3. The errors $\epsilon$ committed by the global weak hypothesis on the local data are calculated by each client and sent to the aggregator.
4. The aggregator exploits the error information to select the best weak hypothesis $c$, adds it to the global strong hypothesis and sends the calculated AdaBoost coefficient $\alpha$ to the collaborators.

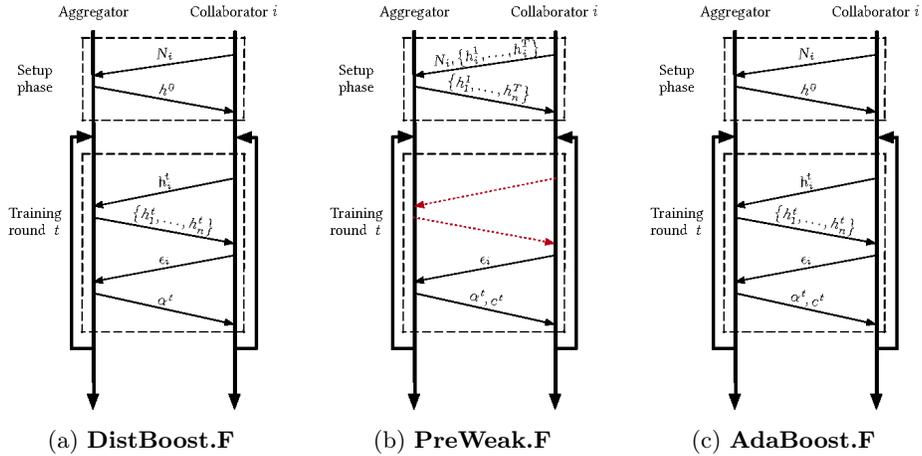(a) **DistBoost.F**          (b) **PreWeak.F**          (c) **AdaBoost.F**

Fig. 1: **The three protocols implied by DistBoost.F, PreWeak.F, and AdaBoost.F**. $N$ is the dataset size, $T$ is the number of training rounds, $h$ the weak hypothesis, $\epsilon$ the classification error, $\alpha$ the AdaBoost coefficient. The subscript $i \in [1, n]$ indices the collaborators and the superscript $t$ the training rounds (with 0 standing for an untrained weak hypothesis). $c \in [1, n]$ is the index of the best weak hypothesis in the hypothesis space. The red dotted line in PreWeak.F indicates the absence of communication.

Note that $N$ is needed to adequately weight the errors committed by the global weak hypothesis on the local data, thus allowing to compute $\alpha$ correctly.

Figure 1 depicts the protocol specialisations for the three algorithms described in [18]. They are similar once abstracted from their low-level details. While step 1 is inherently a setup step, steps 2-4 are repeated cyclically by DistBoost.F and AdaBoost.F. PreWeak.F instead fuses steps 1 and 2 at setup time, receiving from each collaborator $T$ instances of already trained weak hypotheses (one for each training round) and broadcasting $n \times T$ models to the federation. Then, each federated round $t$ loops only on steps 3 and 4 due to the different *hypothesis space* the algorithms explore. While DistBoost.F and AdaBoost.F create a weak hypothesis during each federated round, PreWeak.F creates the whole hypothesis space during step 2 and then searches for the best solution in it.

All three algorithms produce the same strong hypothesis and AdaBoost model, but they differ in the selection of the best weak hypothesis at each round:

 – DistBoost.F uses a committee of weak hypotheses;
 – PreWeak.F uses the weak hypotheses from a fully trained AdaBoost model;
 – AdaBoost.F uses the best weak hypothesis trained in the current round.

The generic model-agnostic federated protocol is more complex than the standard FL one. It requires one more communication for each round and the exchange of complex objects across the network (the weak hypotheses), impacting
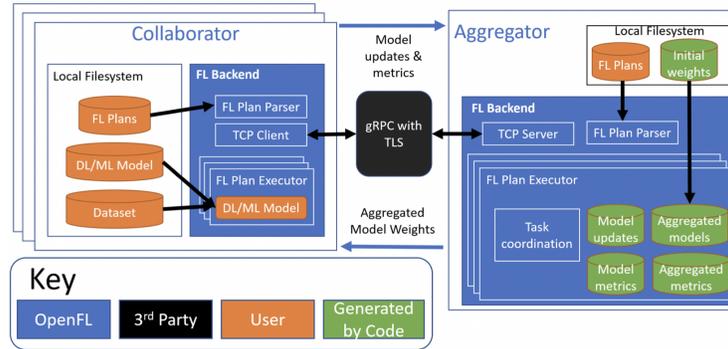
Fig. 2: OpenFL architecture from [5]. The proposed extension targets only the inner components (coloured in blue).

performance. Note that each arrow going from collaborator $i$ to the aggregator in Fig. 1 implies a synchronisation barrier among all the collaborators in the federation. Increasing the number of global synchronisation points reduces concurrency and increases the sensitivity to stragglers. It is worth noting that once an FL framework can handle the common protocol structure, implementing any of the three algorithms requires the same effort. For this study, we implemented AdaBoost.F for two main reasons. First, its protocol covers the whole set of messages (like DistBoost.F), making it computationally more interesting to analyse than PreWeak.F. Besides, AdaBoost.F achieves the best learning results out of the three, also when data is heavily non-IID across the collaborators.

## 4    MAFL Architecture

Redesigning OpenFL comprises two main goals: allowing more flexible protocol management and making the whole infrastructure model agnostic. During this process, we aimed to make the changes the least invasive and respect the original design principles whenever possible (see Fig.2).

### 4.1    The Plan Generalization

The *Plan* guides the software components' run time. It is a YAML file containing all the directives handling the FL learning task, such as which software components to use, where to save the produced models, how many rounds to train, which tasks compose a federated round, and so on. The original OpenFL Plan is rather primitive in its functions. It is not entirely customisable by the user, and many of its fields are overwritten at run time with the default values. Due to its unused power, the parsing of the plan file has been extended and empowered, making it capable of handling new types of tasks, along with a higher range of arguments (and also making it evaluate *every* parameter in the file).

The new model-agnostic workflow can be triggered by specifying the `nn: False` argument under the `Aggregator` and `Collaborator` keywords. The specific steps of the protocol can then be specified in the `tasks` section. In the Intel® OpenFL framework, there are only three possible tasks:

– `aggregated_model_validation`: test set validation of aggregated model;
– `train`: local training of the model;
– `locally_tuned_model_validation`: test set validation of local model.

The three tasks are executed cyclically, with the Aggregator broadcasting the aggregated model before the first task and gathering the local models after the training step. In MAFL, the tasks vocabulary comprises three additional tasks:

– `weak_learners_validate`: test set validation of the weak learners;
– `adaboost_update`: update of the global parameters of AdaBoost.F on the Collaborators and the ensemble model on the Aggregator;
– `adaboost_validate`: local test set validation of the aggregated AdaBoost.F model.

The `weak_learners_validate` task is similar to `aggregated_model_validation`. However, it returns additional information for AdaBoost.F, such as which samples are correctly predicted/mispredicted and the norm of the samples' weights.

The extended set of tasks allows users to use new FL algorithms, such as AdaBoost.F. Additionally, if the `adaboost_update` task is omitted, it is possible to obtain a simple *Federated Bagging* behaviour. Switching behaviour requires small actions other than changing the Plan; however, both functionalities are documented with tutorials in the code repository.

### 4.2   Expanded Communication Protocol

New messages have been implemented into the original *communication protocol*, allowing the exchange of values other than ML models and performance metrics since AdaBoost.F relies on exchanging locally calculated parameters. Furthermore, Intel® OpenFL only implements two synchronisation points in its original workflow: one at the end of the federation round and one when the Collaborator asks the Aggregator for the aggregated model. These synchronisation points are hard-coded into the software and cannot be generalised for other uses.

For the AdaBoost.F workflow, a more general synchronisation point is needed: not two consecutive steps can be executed before each Collaborator has concluded the previous one. Thus a new `synch` message has been added to the *gRPC* protocol. The working mechanism of this synchronisation point is straightforward: the collaborators ask for a `synch` at the end of each task, and if not all collaborators have finished the current task, it is put to sleep; otherwise, it is allowed to continue to the next task. This solution, even if not the most efficient, respects the Intel® OpenFL internal synchronisation mechanisms and thus does not require any different structure or new dependency.

### 4.3   Core Classes Extension

The following core classes of the framework have been modified to allow the standard and model-agnostic workflows to coexist (see Fig. 2 for an overview).

The `Collaborator` class can now offer different behaviours according to the ML model used in the computation. Suppose the Plan specifies that the training will not involve DNNs. In that case, the Collaborator will actively keep track of the parameters necessary to the AdaBoost.F algorithm, like the mispredicted examples, the weight associated with each data sample, and the weighted error committed by the models. Additionally, the handling of the internal database used for storage will change behaviour, changing tags and names associated with the entries to make possible finer requests to it.

The `Aggregator` can now generate any ML models (instead of only DNNs weights), handle aggregation functions instantiated dynamically from the plan file, and handle the synchronisation needed at the end of each step. New methods allow the Aggregator to query the internal database more finely, thus allowing it to read and write ML models with the same tags and name as the Collaborator.

`TensorDB`, the internal class used for storage, has been modified to accommodate the new behaviours described above. This class implements a simple *Pandas* data frame responsible for all model storage and retrieving done by the Aggregator and Collaborators. Furthermore, its `clean_up` method has been revised, making it possible to maintain a fixed amount of data in memory. This fix has an important effect on the computational performance since the query time to this object is directly proportional to the amount of data it contains.

Finally, the more high-level and interactive classes, namely `Director` and `Envoy`, and the serialization library have been updated to work correctly with the new underlying code base. These software components are supposed to be long-lived: they should constantly be running on the server and clients' hosts. When a new experiment starts, they will instantiate the necessary `Aggregator` and `Collaborators` objects with the parameters for the specified workflow.

This effort results in a model-agnostic FL framework that supports the standard DNNs-based FL workflow and the new AdaBoost.F algorithm. Using the software in one mode or another does not require any additional programming effort from the user: a few simple configuration instructions are enough. Additionally, the installation procedure has been updated to incorporate all new module dependencies of the software. Finally, a complete set of tutorials has been added to the repository: this way, it should be easy for any developer to get started with this experimental software.

## 5   Evaluation

The complete set of tutorials replicating the experiments from [18] are used to assess MAFL's correctness and efficiency. We run them on a cloud and HPC infrastructure, both x86-64 based, and Monte Cimone, the first RISC-V based HPC system; however, MAFL runs also on ARM-v8 and Power systems.

### 5.1  Performance Optimizations

Using weak learners instead of DNNs drastically reduces the computational load. As an example, [1] reports 18.5 vs 419.3 seconds to train a 10-leaves decision tree or a DNN model, respectively, on the PRAISE training set (with comparable prediction performance). Moreover, AdaBoost.F requires one additional communication phase per round. This exacerbates the impact of time spent in communication and synchronisation on the overall system performance. To reduce this impact, we propose and evaluate different optimisations to reduce this overhead. Applying all proposed optimisations, we achieve a 5.5x speedup on a representative FL task (see Fig. 3). As a baseline workload, we train a 10-leaves decision tree on the Adult dataset over 100 rounds using 9 nodes (1 aggregator plus 8 collaborators). We use physical machines to obtain stable and reliable computing times, as execution times on bare-metal nodes are more deterministic than cloud infrastructures. Each HPC node is equipped with two 18-core Intel® Xeon E5-2697 v4 @2.30 GHz and 126 GB of RAM. A 100Gb/s Intel® Omni-Path network interface (in IPoFabric mode) is used as interconnection network. Reported times are average of five runs $\pm$ the 95% CI.

We start by measuring the execution time given by the baseline: 484.13±15.80 seconds. The first optimisation is to adapt the buffer sizes used by gRPC to accommodate larger models and avoid resizing operations. Increasing the buffer from 2MB to 32MB using decision trees reduced the execution time to 477.0±17.5 seconds, an improvement of $\sim$ 1.5%. While this seems small, the larger the models, the bigger the impact of this optimisation. The second optimisation is the choice of the serialisation framework: by using `Cloudpickle`, we reduce the execution time to 471.4±6.1 seconds, an improvement of $\sim$2.6%. Next, we examine `TensorDB`, which grows linearly in the number of federated rounds, thus slowing down access time linearly. We modified the `TensorDB` to store only the essential information of the last two federation rounds: this results in a stable memory occupation and access time. With this change, the execution time drops to 414.8±0.9 seconds, an improvement of $\sim$14.4% over the baseline.

Lastly, two `sleep` are present in the MAFL code: one for the end-round synchronisation and another for the `synch` general synchronisation point, fixed respectively at 10 and 1 seconds. Both have been lowered to 0.01 seconds since we assessed empirically that this is the lowest waiting time still improving the global execution time. This choice has also been made possible due to the computational infrastructures exploited in this work; it may not be suitable for wide-scale implementations in which servers and clients are geographically distant or compute and energy-constrained. With this sleep calibration, we obtained a global execution time of 250.8±9.6 seconds, a $\sim$48.2% less than the baseline. Overall, with all the optimisations applied together, we can achieve a final mean execution time of 88.6± seconds, i.e. a 5.46x speedup over the baseline.

### 5.2  Correctness

We replicate the experiments from [18] and compare the ML results. These experiments involve ten different datasets: `adult`, `forestcover`, `kr-vs-kp`, `splice`,
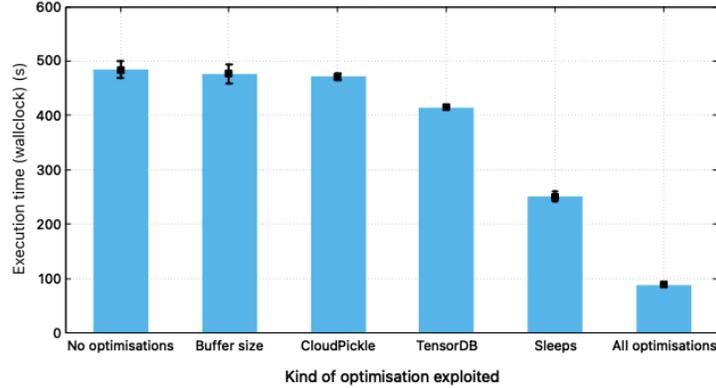
Fig. 3: Ablation study of the proposed software optimisations; the 95% CI has been obtained over five executions.

Table 1: Mean F1 scores ± standard deviation over 5 runs.

| Dataset | Classes | Reference | MAFL |
|---------|---------|-----------|------|
| Adult | 2 | $85.58 \pm 0.06$ | $85.60 \pm 0.05$ |
| ForestCover | 2 | $83.67 \pm 0.21$ | $83.94 \pm 0.14$ |
| Kr-vs-kp | 2 | $99.38 \pm 0.29$ | $99.50 \pm 0.21$ |
| Splice | 3 | $95.61 \pm 0.62$ | $96.97 \pm 0.65$ |
| Vehicle | 4 | $72.94 \pm 3.40$ | $80.04 \pm 3.30$ |
| Segmentation | 7 | $86.07 \pm 2.86$ | $85.58 \pm 0.06$ |
| Sat | 8 | $83.52 \pm 0.58$ | $84.89 \pm 0.57$ |
| Pendigits | 10 | $93.21 \pm 0.80$ | $92.06 \pm 0.44$ |
| Vowel | 11 | $79.80 \pm 1.47$ | $79.34 \pm 3.31$ |
| Letter | 26 | $68.32 \pm 1.63$ | $71.13 \pm 2.02$ |

`vehicle`, `segmentation`, `sat`, `pendigits`, `vowel`, and `letter`. These are standard ML datasets targeting classification tasks, both binary (`adult`, `forestcover`, `kr-vs-kp`) and multi-class (all the others), with a varying number of features (from the 14 of `adult` up to the 61 of `splice`), and a different number of samples (from the 846 of `vehicle` up to the 495.141 of `forestcover`). Each training set has been split in an IID way across all the Collaborators, while the testing has been done on the entire test set. A simple Decision Tree from SciKit-Learn with ten leaves is used as a weak learner; instead, the AdaBoost class has been created manually. We set the number of federated rounds to 300 and use 10 nodes: 1 aggregator plus 9 collaborators. We note that these optimizations can also benefit the original OpenFL.

Table 1 reports each dataset's reference and calculated F1 scores (mean value ± the standard deviation over five runs). The values reported are fully compatible with the results reported in the original study, thus assessing the correctness

of the implementation. In particular, it can be observed that the standard deviation intervals are particularly high for the `vehicle`, `segmentation`, and `vowel`. This fact can be due to the small size of the training set of these datasets, respectively 677, 209, and 792 samples, which, when split up across ten Collaborators, results in an even smaller quantity of data per client: this can thus determine the creation of low-performance weak learners. Furthermore, also `letter` reported a high standard deviation: this could be due to the difference between the classification capabilities of the employed weak learner (a 10-leaves Decision Tree) compared to the high number of labels present in this dataset (26 classes), thus making it hard to obtain high-performance weak learners.

The mean F1 score curve for each dataset can be observed in Figure 4a. As can be seen, after an initial dip in performance, almost each learning curve continues to grow monotonically to higher values. This fact is expected since the AdaBoost.F is supposed to improve its classification performance with more weak learners. It has to be observed that, at each federated round, a new weak learner will be added to the aggregated model: the AdaBoost.F grows linearly in size with the number of federated rounds. This characteristic of the algorithm has many consequences, like the increasingly longer time needed for inference and for moving the aggregated model over the network. From Figure 4a, we can observe that, in the vast majority of cases, a few tens of federated rounds are more than enough to obtain a decent level of F1 scores; this is interesting since it is possible to obtain a small and efficient AdaBoost.F model in little training effort. Instead, for the more complex datasets like `letter` and `vowel`, we can observe that it is possible to obtain better performance with longer training efforts. This means that is possible to use AdaBoost.F to produce bigger and heavier models at need, according to the desired performance and inference complexity.

### 5.3   Flexibility

To demonstrate the model-agnostic property of MAFL, we choose the `vowel` dataset and train different ML model types on it. In particular, one representative ML model has been chosen from each multi-label classifier family available on SciKit-Learn: Extremely Randomized Tree (Trees), Ridge Linear Regression (Linear models), Multi-Layer Perceptron (Neural Networks), K-Nearest Neighbors (Neighbors), Gaussian Naive Bayes (Naive Bayes), and simple 10-leaves Decision Trees as baselines. Fig. 4b summarises the F1 curves for the different ML models used as weak learners. Each model has been used out-of-the-box, without hyper-parameter tuning using the default parameters set by SciKit-Learn v1.1.2. All ML models work straightforwardly in the proposed software without needing to code anything manually: it is sufficient to replace the class name in the experiment file. This proves the ease with which data scientists can leverage MAFL to experiment with different model types.
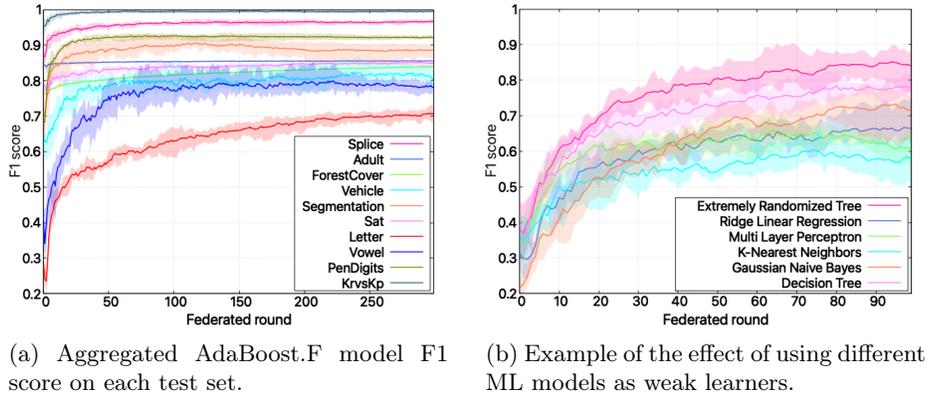
(a) Aggregated AdaBoost.F model F1 score on each test set.

(b) Example of the effect of using different ML models as weak learners.

Fig. 4: ML properties of MAFL

## 5.4   Scalability Analysis

We perform this scalability study using the HPC nodes from Sec. 5.1 and Monte Cimone [2], the first available HPC-class cluster based on RISC-V processors. It comprises eight computing nodes equipped with a U740 SoC from SiFive integrating four U74 RV64GCB cores @ 1.2 GHz, 16GB RAM and a 1 Gb/s interconnection network.

We select the `forestcover` dataset for running these experiments, being the largest dataset used in this study, split into a 485K training samples and 10K testing samples. The weak learner is the same 10-leaves SciKit-Learn Decision Tree from Sec.5.2. We lowered the number of federated training rounds to 100 since they are enough to provide an acceptable and stable result (10 on the RISC-V system due to the longer computational times required). Different federations have been tested, varying numbers of Collaborators from 2 to 64 by powers of 2. We went no further since OpenFL is designed to suit a cross-silo FL scenario, which means a few tens of clients. We investigated two different scenarios: *strong scaling*, where we increase the collaborators while keeping the same problem size by spitting the dataset samples in uniform chunks across collaborators; and *weak scaling*, where we scale the problem size with the number collaborators by assigning each collaborator the entire dataset. In both cases, the baseline reference time is the time taken by a federation comprising the aggregator and a single collaborator. We report the mean over 5 runs for each experiment.

Fig. 5 shows the strong and weak scaling properties of MAFL. The RISC-V plot stops at 7 because we have just 8 nodes in the cluster, and we want to avoid sharing a node between the aggregator and collaborator to maintain the same experiment system setting. In the strong scaling scenario, the software does not scale efficiently beyond 8 HPC nodes, as the execution becomes increasingly communication-bound. The same also affects the weak scaling. Nevertheless, the degradation is sublinear (each point on the x/axis doubles the number of nodes). This is important because the main benefit in the FL scenario is the additional

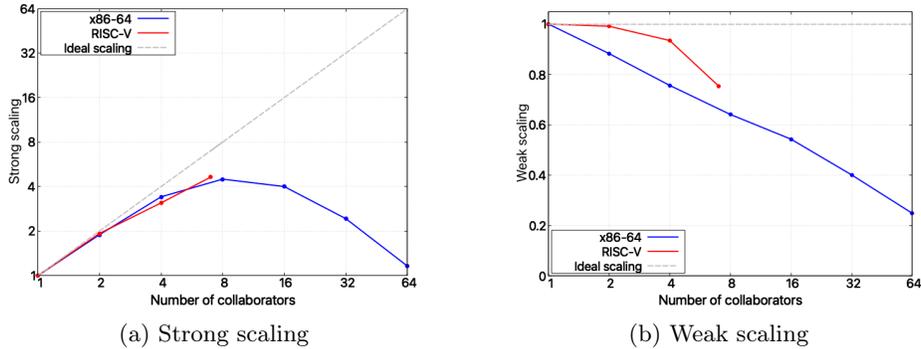(a) Strong scaling  (b) Weak scaling

Fig. 5: Strong and weak scaling properties of MAFL.

training data brought in by each contributor node. The RISC-V cluster exhibits better strong scalability when comparing the two clusters. This is justified by the slower compute speed of the RISC-V cores leading to higher training times, making the execution more compute-bound, especially for a low number of nodes. The weak scalability on the RISC-V cluster suffers from the lower network speed. Since real-world cross-silo federations rarely count more than a dozen participants, it can be assessed that MAFL is suitable for experimenting with such real-world scenarios.

## 6 Discussion

The implementation experience of MAFL and the subsequent experimentation made it evident that current FL frameworks are not designed to be as flexible as the current research environment needs them to be. The fact that the standard workflow of OpenFL was not customisable in any possible way without modifying the code and that the serialisation structure is DNN-specific led the authors to the idea that a new, workflow-based FL framework is needed. Such a framework should not implement a fixed workflow but allow the user to express any number of workflow steps, entities, the relations between them, and the objects that must be exchanged. This property implies the generalisation of the serialisation infrastructure, which cannot be limited to tensors only. Such an approach would lead to a much more straightforward implementation of newer and experimental approaches to FL, both from the architectural and ML perspective.

Furthermore, the use of asynchronous communication can help better manage the concurrent architecture of the federation. These systems are usually slowed down by stragglers that, since the whole system is supposed to wait for them, will slow down the entire computation. In our experience implementing MAFL, a significant part of the scalability issues is determined by the waiting time between the different collaborators taking part in the training. While such an approach would improve the scalability performance of any FL framework, it also underlies the investigation of how to simultaneously handle newer and older updates. This

capability would improve the computational performance of gradient and non-gradient-based systems: the relative aggregation algorithms must be revised to accommodate this new logic. This matter is not trivial and deserves research interest. Lastly, due to the possibility of exploiting less computationally requiring models, MAFL can easily be used to implement FL on low-power devices, such as systems based on the new RISC-V.

## 7    Conclusions

A model-agnostic modified version of Intel$^{\circledR}$ OpenFL implementing the Ad-aBoost.F federated boosting algorithm, named MAFL, has been proposed. The proposed software has been proven to implement the AdaBoost.F algorithm correctly and can scale sufficiently to experiment efficiently with small cross-silo federations. MAFL is open-source, freely available online, easily installable, and has a complete set of already implemented examples. To our knowledge, MAFL is the first FL framework to implement a model-agnostic, non-gradient-based algorithm. This effort will allow researchers to experiment with this new conception of FL more freely, pushing the concept of model-agnostic FL even further. Furthermore, this work aims to contribute directly to the RISC-V community, enabling FL research on this innovative platform.

## References

1. Arfat, Y., Mittone, G., Colonnelli, I., D'Ascenzo, F., Esposito, R., Aldinucci, M.: Pooling critical datasets with federated learning. In: IEEE PDP (2023)
2. Bartolini, A., Ficarelli, F., Parisi, E., Beneventi, F., Barchi, F., Gregori, D., et al.: Monte cimone: Paving the road for the first generation of risc-v high-performance computers. In: IEEE SOCC. pp. 1–6 (2022)
3. Beltrán, E.T.M., Pérez, M.Q., Sánchez, P.M.S., Bernal, S.L., Bovet, G., Pérez, M.G., et al.: Decentralized federated learning: Fundamentals, state-of-the-art, frameworks, trends, and challenges. arXiv preprint arXiv:2211.08413 (2022)
4. Beutel, D.J., Topal, T., Mathur, A., Qiu, X., Parcollet, T., de Gusmão, P.P., et al.: Flower: A friendly federated learning research framework. arXiv preprint arXiv:2007.14390 (2020)
5. Foley, P., Sheller, M.J., Edwards, B., Pati, S., Riviera, W., Sharma, M., et al.: Openfl: the open federated learning library. Phys. Med. Biol. **67**(21), 214001 (2022)
6. Freund, Y., Schapire, R.E.: A decision-theoretic generalization of on-line learning and an application to boosting. J. Comput. Syst. Sci. **55**(1), 119–139 (1997)

7. He, C., Li, S., So, J., Zhang, M., Wang, H., Wang, X., et al.: Fedml: A research library and benchmark for federated machine learning. arXiv preprint arXiv:2007.13518 (July 2020)

8. Holzinger, A., Langs, G., Denk, H., Zatloukal, K., Müller, H.: Causability and explainability of artificial intelligence in medicine. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery **9**(4), 1312 (2019)

9. Kairouz, P., McMahan, H.B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A.N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al.: Advances and open problems in federated learning. Found. Trends Mach. Learn. **14**(1-2), 1–210 (2021)

10. Kleanthous, C., Chatzis, S.: Gated mixture variational autoencoders for value added tax audit case selection. Knowl. Based Syst. **188**, 105048 (2020)

11. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. Commun. ACM **60**(6), 84–90 (2017)

12. Liu, Y., Fan, T., Chen, T., Xu, Q., Yang, Q.: Fate: An industrial grade platform for collaborative learning with data protection. J. Mach. Learn. Res. **22**(1), 10320–10325 (2021)

13. Ludwig, H., Baracaldo, N., Thomas, G., Zhou, Y., Anwar, A., Rajamoni, S., et al.: IBM federated learning: an enterprise framework white paper v0. 1. arXiv preprint arXiv:2007.10987 (2020)

14. Lyu, L., Yu, H., Ma, X., Chen, C., Sun, L., Zhao, J., et al.: Privacy and robustness in federated learning: Attacks and defenses. IEEE Trans. Neural. Netw. Learn. Syst. pp. 1–21 (2022)

15. McMahan, B., Moore, E., Ramage, D., Hampson, S., Agüera y Arcas, B.: Communication-efficient learning of deep networks from decentralized data. In: Proceedings of the 20th International Conference on Artificial Intelligence and Statistics AISTATS. vol. 54, pp. 1273–1282. PMLR, Fort Lauderdale, FL, USA (2017)

16. Meese, C., Chen, H., Asif, S.A., Li, W., Shen, C.C., Nejad, M.: BFRT: Blockchained federated learning for real-time traffic flow prediction. In: IEEE CCGrid. pp. 317–326 (2022)

17. O'Mahony, N., Campbell, S., Carvalho, A., Harapanahalli, S., Hernandez, G.V., Krpalkova, L., et al.: Deep learning vs. traditional computer vision. In: SAI. vol. 943, pp. 128–144. Springer, Cham (2019)

18. Polato, M., Esposito, R., Aldinucci, M.: Boosting the federation: Cross-silo federated learning without gradient descent. In: IEEE IJCNN). pp. 1–10 (2022)

19. Riviera, W., Menegaz, G., Boscolo Galazzo, I.: FeLebrities: a user-centric assessment of federated learning frameworks. TechRxiv (2022)

20. Roth, H.R., Cheng, Y., Wen, Y., Yang, I., Xu, Z., Hsieh, Y.T., et al.: Nvidia flare: Federated learning from simulation to real-world. arXiv preprint arXiv:2210.13291 (2022)

21. Sotthiwat, E., Zhen, L., Li, Z., Zhang, C.: Partially encrypted multi-party computation for federated learning. In: IEEE CCGrid. pp. 828–835 (2021)

22. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: NeurIPS. pp. 3104–3112 (2014)

23. Warnat-Herresthal, S., Schultze, H., Shastry, K.L., Manamohan, S., Mukherjee, S., Garg, V., et al.: Swarm learning for decentralized and confidential clinical machine learning. Nature **594**(7862), 265–270 (2021)

24. Zhavoronkov, A., Ivanenkov, Y.A., Aliper, A., Veselov, M.S., Aladinskiy, V.A., Aladinskaya, A.V., et al.: Deep learning enables rapid identification of potent ddr1 kinase inhibitors. Nature biotechnology **37**(9), 1038–1040 (2019)