

---

# iFedAvg – Interpretable Data-Interoperability for Federated Learning

---

**David Roschewitz**  
ETH Zürich  
Switzerland  
david.roschewitz@inf.ethz.ch

**Mary-Anne Hartley**  
EPFL, Lausanne  
Switzerland  
mary-anne.hartley@epfl.ch

**Luca Corinzia**  
ETH Zürich  
Switzerland  
luca.corinzia@inf.ethz.ch

**Martin Jaggi**  
EPFL, Lausanne  
Switzerland  
martin.jaggi@epfl.ch

## Abstract

Recently, the ever-growing demand for privacy-oriented machine learning has motivated researchers to develop federated and decentralized learning techniques, allowing individual clients to train models collaboratively without disclosing their private datasets. However, widespread adoption has been limited in domains relying on high levels of user trust, where assessment of data compatibility is essential. In this work, we define and address low interoperability induced by underlying client data inconsistencies in federated learning for tabular data. The proposed method, iFedAvg, builds on federated averaging adding local element-wise affine layers to allow for a personalized and granular understanding of the collaborative learning process. Thus, enabling the detection of outlier datasets in the federation and also learning the compensation for local data distribution shifts without sharing any original data. We evaluate iFedAvg using several public benchmarks and a previously unstudied collection of real-world datasets from the 2014 - 2016 West African Ebola epidemic, jointly forming the largest such dataset in the world. In all evaluations, iFedAvg achieves competitive average performance with negligible overhead. It additionally shows substantial improvement on outlier clients, highlighting increased robustness to individual dataset shifts. Most importantly, our method provides valuable client-specific insights at a fine-grained level to guide interoperable federated learning.

## 1 Introduction

Institutions with sensitive data, such as hospitals, cannot typically share patient data due to privacy regulations. However, solely relying on in-house data can lead to models with poor generalization due to the limited, and potentially biased, input data. Federated learning (FL) partially addresses this issue by enabling various clients to contribute and benefit from a collaborative learning process without revealing their underlying data. In practice though, individual clients can benefit from the collaborative training only if their data is compatible with that of other participating institutions. This can lead to situations where the client data is not interoperable, and where joining the federated learning process yields no benefits or even has a detrimental effect. Furthermore, a lack of transparency of the federated learning process impairs the trust of the federation, limiting its adoption by more institutions.

Thus, there is a clear need to address the interoperability of current federated learning approaches. Such methods should not only detect potential data shifts and automatically correct them but, more importantly, they should also be easily interpretable for stakeholders to visually assess the suitability of the collaboration.

We study these challenges and propose a novel method for a setting inspired by a real-world medical dataset collected during the 2014-16 Ebola epidemic. The data was collected at different treatment centres, by various organizations in multiple countries generating an inherently heterogeneous dataset. In practice, no single agent would have access to the data of other agents, necessitating a federated learning approach. Further details are outlined in Section 3. Our method is tailored to handle tabular datasets which are abundant in practice and for which feature-shifts are intuitive to understand.

The fundamental approach that we consider is to learn a personalized data transformation for each client during the federated training procedure. This transformation can be viewed as a local re-normalization or embedding that makes clients more interoperable without ever exchanging data. Through deliberate design, our method attains unparalleled transparency, allowing fine-grained interpretation of the learned shifts for each feature of each client relative to their peers in the federation. This ability to compare data shifts across clients means that interoperability can be easily assessed. For example, clients collecting data in pediatric, adult or geriatric medicine would not only have the "age" feature highlighted as responsible for their "outlier" status, but the directionality and relative magnitude of their outlier shifts can also be assessed.

Our main contributions are the following:

1. We propose a novel framework, iFedAvg which detects and corrects interoperability issues in federated learning on tabular datasets.
2. We present visualization tools for practitioners to assess the feature-wise compatibility to collaborative learning.
3. Finally, we demonstrate the potential of the proposed method on a previously unstudied collection of data from the West African Ebola epidemic and multiple public benchmarks.

**Outlook** In Section 2 we outline our method, and we present a detailed introduction to the Ebola dataset in Section 3. Subsequently we present the experimental setup and results in Sections 4 and 5. We provide a discussion and conclusions in the last section.

**Related Work** The concept of federated learning was formulated by McMahan et al. [20] as “*collaborative machine learning without centralized training data*” alongside *Federated Averaging* (FedAvg). The method proved efficient in learning a single, global, gradient-based model from many clients’ data in a private fashion. However, this approach is impaired in the realistic setting with statistically heterogeneous client data Zhao et al. [28]. In order to address some of these shortcomings, multiple extensions have been proposed. For instance, FedProx by Li et al. [19] allows for inexact local computations and regularizes client drift using an additive proximal term in the loss function. Karimireddy et al. [17] introduced SCAFFOLD, which explicitly uses control variates to improve convergence. From an optimizer perspective, SGD with server-side momentum has proven effective [11], a result further investigated by Reddi et al. [22] that proposes the use of adaptive optimizers for federated learning such as FedAdam. For a more comprehensive overview, we refer to the excellent review of federated learning by Kairouz et al. [16].

Learning a client-specific *data transformation*, as we propose here, can be seen as a special case of training personalized models for each client as opposed to a single global model. In the area of personalized federated learning (PFL), various approaches have been identified. A relevant selection includes transfer learning, multi-task learning, meta-learning and personalization layers. Transfer and meta-learning focus on tuning an initial model, usually the global federated model, to each client. Techniques here include fine-tuning Yu et al. [27] and Model-Agnostic-Meta-Learning (MAML) Jiang et al. [15]. Fallah et al. [7] devise Per-FedAvg, leveraging a second-order derivative to account for personalization throughout the process. Multi-task learning approaches focus on jointly learning multiple models for a variety of tasks with different levels of similarity. This can be applied to PFL, as in MOCHA [24], or using a Bayesian framework [4]. Furthermore, personalization can be achieved by client-specific layers, as proposed by [2], who show that locally trained output layers are effective for image classification. Furthermore, Deng et al. [5] propose a method, APFL, which actively optimizes a global and a local model, and blends the two models concurrently. The summary overview provided by Kulkarni et al. [18] includes additional approaches for PFL.

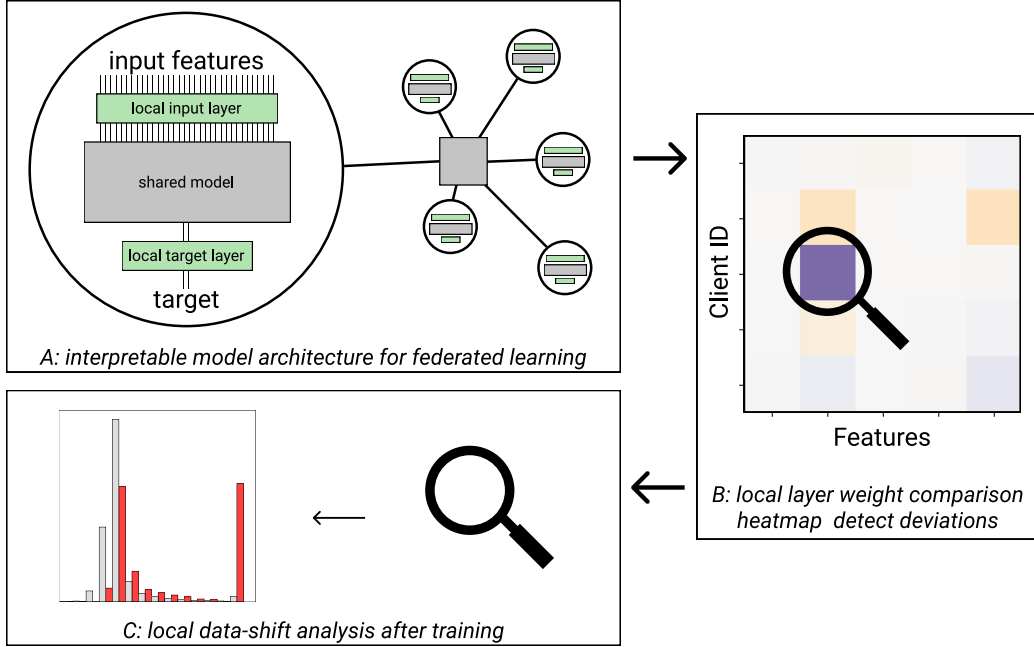


Figure 1: The three phases enabled by iFedAvg. A) a novel model architecture that extracts feature-wise interoperability information deployable in a federated setting B) interpretable outputs allowing practitioners to detect and understand inter-client compatibility issues C) independent private interpretation performed after the training procedure.

Understanding *differences* between clients and interpreting the federated learning process as a whole is, to the best of our knowledge, a problem not previously investigated. Standard neural-network model-interpretation techniques such as LIME [23] are, unfortunately, not immediately applicable to federated learning and thus would not fulfil our objectives. For vertically (i.e., feature-wise) partitioned data and tree-based models, SHAP-values have been investigated by Zheng et al. [29] and Wang [25], but the approach is not applicable to our setting as we focus on a *horizontally* (i.e., sample-wise) partitioned dataset. Furthermore, SHAP values learn feature-wise contributions to the model, rather than learning and compensating for potential data biases. Imakura et al. [12] present an "interpretable non-model sharing collaborative data analysis method as a federated learning system". While retaining privacy, the authors assume having access to a shared public *anchor* dataset and focus on model-interpretability only. Hence the need for a method allowing for interpretability at the process level enabling comparisons between clients.

## 2 Interpretable and data-interoperable federated averaging

Our proposed method, iFedAvg is designed not only as a personalized federated learning algorithm but as a component of a complete framework. Figure 1 illustrates this entire workflow, highlighting how iFedAvg enables each step. The following subsections outline the main aspects of the proposed method.

### 2.1 Architecture

Let us denote an existing neural network model which could be shared among all clients as  $f_{\text{shared}} : \mathbb{R}^D \rightarrow \mathbb{R}^K$ . This model maps an input vector  $\mathbf{x} \in \mathbb{R}^D$  to a target vector  $\mathbf{y} \in \mathbb{R}^K$ . We place no assumptions on the type or complexity of the  $f_{\text{shared}}$ 's network as long as it can be trained using gradient-based optimizers and conforms to the definition above. The objective of our work is to devise an extension, allowing for  $f_{\text{shared}}$  to be more interoperable in a transparent fashion.

iFedAvg introduces personalization layers around the shared neural network,  $f_{\text{shared}}$ . The combined model is then specified as  $f_{\text{out}} \circ f_{\text{shared}} \circ f_{\text{in}}$ , where  $\circ$  indicates a composition. To retain the correct dimensionality, the input and output layers are as follows:  $f_{\text{in}} : \mathbb{R}^D \rightarrow \mathbb{R}^D$  and  $f_{\text{out}} : \mathbb{R}^K \rightarrow \mathbb{R}^K$ .

In order to retain interpretability at the feature- and target-level, we propose the layers to simply be an element-wise learned normalization. We assume that numerical values are standardized for each client. Inspired by traditional standardization we explicitly define  $f_{\text{in}}$ , with bias and weight vectors  $\mathbf{b}_{\text{in}}, \mathbf{w}_{\text{in}} \in \mathbb{R}^D$  as:

$$f_{\text{in}}(\mathbf{x}) = (\mathbf{x} + \mathbf{b}_{\text{in}}) \odot \mathbf{w}_{\text{in}}, \quad (1)$$

where  $\odot$  refers to the element-wise multiplication. This construction does not allow the blending that a traditional *fully-connected* layer permits.  $f_{\text{in}}$  is initialized to be the identity function, namely setting  $\mathbf{b}_{\text{in}} = 0$  and  $\mathbf{w}_{\text{in}} = 1$ . Similarly, we can define  $f_{\text{out}}(\mathbf{y}) = (\mathbf{y} + \mathbf{b}_{\text{out}}) \odot \mathbf{w}_{\text{out}}$ , with  $\mathbf{b}_{\text{out}}, \mathbf{w}_{\text{out}} \in \mathbb{R}^K$ . Furthermore, with only additional  $2D + 2K$  parameters, the memory cost is negligible.

At first glance, our approach appears analogous to the *personalization layers* proposed by Arivazhagan et al. [2]. iFedAvg differs however, in the purposeful placement and heavily restricted design of the layers which enable actionable insights to be extracted. Furthermore, this study does not solely analyze the personalization properties of the method, but focuses on revealing how the federated learning process deals with biases in underlying data as well as offering a means of compensating for them.

## 2.2 Training

Training iFedAvg does not differ substantially from FedAvg. At each round, each client performs a local update using stochastic gradient descent (SGD) and locally retains the updates on  $f_{\text{in}}$  and  $f_{\text{out}}$ . Updates for all weights of  $f_{\text{shared}}$  are disclosed to the server, which performs the standard FedAvg aggregation step and broadcasts the new shared weights to all clients.

This procedure is significantly more efficient than other related personalization methods such as APFL or PER-FedAvg. From a performance perspective, neither a second backpropagation or a Hessian, respectively, need to be computed. A single iteration of SGD is able to update the entire *combined* model at once. Furthermore, only a single copy of the network needs to be stored, which reduces the storage demands on each client, compared to APFL.

Intuitively, these layers allow each client to learn a feature-shift as well as target-shift, with respect to the federation. As the central block of the neural network is shared, and therefore identical for each client, the personalized layers can be seen as learning the necessary transformation from the underlying private data to the model of the federation.

## 2.3 Interpretable outputs

The cornerstone of the interpretability of iFedAvg is the personalized layer design. The layers learn the local shifts necessary to be able to utilize the shared model block of the federation. Each client can, therefore, adjust their *combined* model in a very restricted sense. This restriction-by-design means that each individual value in  $f_{\text{in}}$  and  $f_{\text{out}}$  is directly interpretable.

For instance, following the example above, if the age of patients varies from client to client, clients are, at first, indistinguishable after standardization. With our method, however, the necessary personalized age shift can be learned seamlessly throughout the process and can be learned to correctly predict the diagnosis. The magnitude and direction of this shift is, by design, tied directly to the input feature, and provides unparalleled insights about each participant of the federation without sharing any original data.

A crucial distinction must be made between our method and comparing the underlying data distribution a-priori. Exchanging the datasets' summary statistics can be a useful way to diagnose interoperability, but is limited in three critical ways. First, every feature must be analyzed for every client - in practice the needed effort might render this impossible. Second, revealing summary statistics might still be problematic from a data privacy perspective - personalized model weights are more secure in that regard. Third, iFedAvg only learns *necessary* shifts which aid model performance, eliminating the necessity to investigate potential shifts for insignificant features.



The final step of our proposed architecture is a communication of the learned shifts to stakeholders. Large compensations could be an indicator of critical differences in data collection, calibration, missingness-not-at-random or otherwise. A decision on whether this means a client’s dataset is incompatible or, with the shift-adjustment, is suitable is left to domain experts. *iFedAvg* strives to be the best tool to make an informed decision for which clients and features to include in their federation. Significant deviations could then also allow targeted queries between users to better guide collaboration.

### 3 Federated Ebola dataset

The 2014-16 Ebola Virus Disease (EVD) epidemic in West Africa revealed the devastating consequences of inadequate data sharing during public health emergencies. Delays and poorly compatible datasets were held directly responsible for the slow response to the outbreak, ultimately exacerbating the epidemic Georgetown University Medical Center [8]. In response, the Infectious Disease Data Observatory (IDDO) was established to collate and align the fragmented and poorly interoperable datasets into the largest central repository of Ebola data in the world (i.e., the Ebola Data Platform, EDP) to facilitate coordinated research [13].

The commendable EDP initiative necessitated a laborious process of acquiring ethical approval as well as devising a common data sharing protocol. While this was ultimately highly successful, it took several years before being made available to researchers in 2019, well after it would have been useful in forming an evidence-based response to the crisis. Thus, the EDP is an ideal real-world use-case for exploring the potential of an interoperability-adjusted federated learning approach with the end goal of enabling secure real-time model sharing in rapidly evolving public health emergencies which suffer from poor response coordination.

The EDP comprises tabular clinical data on 13552 anonymized patients treated at 16 Ebola Treatment Centres (ETCs) between January 2014 and December 2015. The ETCs were scattered across the three main affected countries of Sierra Leone, Guinea and Liberia which differ in language, geography, epidemiology, demographics and treatment protocols. Thus, there is high risk for bias ranging from natural variation (e.g., malaria prevalence) to measurement errors (e.g., different tools/protocols to quantify viral load) and misattribution (e.g., mislabelling or poor standardization) and variable missingness. The collected data includes both categorical and continuous features such as demographic details (e.g., age, sex, location), clinical signs and symptoms (e.g., fever, coughing, headache), laboratory values (e.g., Ebola test results and quantitation of viral load) and outcomes for each patient (e.g., death vs recovery). Research on such data is often focused on making diagnostic and prognostic models to better allocate limited resources to the most critical patients and improve early case identification [9, 10, 3, 14].

However, these studies were performed on single datasets where statistical power is diluted in the small numbers of included patients.

As a proof of concept, we replicate these studies, by learning diagnostic and prognostic classification tasks (respectively, EVD negative vs EVD positive and survival vs death in EVD positive cases). Here, the ETC site represents a client with its locally collected dataset and we thus explore the potential of *iFedAvg* to create a robust personalized predictive model while detecting and compensating for known interoperability issues between sites.

The ethical framework and anonymization protocols are published on the IDDO website [13]. We provide further details such as sample sizes and class imbalance in Appendix A.5.

### 4 Experimental design

**Evaluation.** To fairly evaluate *iFedAvg*, we compare it to several state-of-the-art methods which have similar aims, and explore their limitations in a standardized realistic setting. Specifically, each client can choose to not partake in the federation, and simply train a local model. Likewise, vanilla federated averaging, *FedAvg* is a valuable benchmark, as well as a more sophisticated personalized federated algorithm such as Adaptive Personalized Federated Learning, APFL, proposed by Deng et al. [5]. An interesting non-personalized baseline is a single model trained on a centralized concatenation of all clients’ datasets (called in the following the Centralized method). While this is not a realistic

scenario, it is currently the yardstick of many large scale data sharing efforts such as IDDO, which makes it a particularly appropriate benchmark for the Ebola dataset.

Four datasets will serve as the foundation of our experiments:

- **Ebola Prognosis:** Predicting the survival of EVD-positive patients; [13].
- **Ebola Diagnosis:** Predicting whether a patient triaged as "suspect" has EVD; [13].
- **Vehicle Classification with Sensor Network (VSN):** Classifying the type of vehicle based on a network of 23 acoustic and seismic sensors; [6].
- **Human Activity Recognition (HAR):** Classifying the type of activity performed by 30 human subjects according to readings from body sensors; [1].

Every method outlined above is trained on each dataset, with each client retaining 33% or minimally 100 samples of its local data as a hold-out test. Performance metrics are computed locally on this hold-out set, retaining an understanding of personalized performance.

**Preprocessing & missing values.** The experimental setting assumes semantic interoperability, meaning all client datasets were aligned in feature nomenclature. All numerical values are standardized to mean 0 and standard deviation 1. In order to include missingness as an assessed feature of interoperability (i.e., to detect whether the bias is due to non-random missing values), missing continuous values were filled with 0s, binary features with the value 0.5 following standardization.

**Architecture.** For every method, an MLP model with identical architecture is used, only modifying the training regime for each algorithm. In the main experimental results we only enable  $f_{in}$  to investigate feature-shifts. For results and a thorough discussion of iFedAvg with  $f_{out}$ , we refer to Appendix A.4.

**Training.** Each round, every client performs one training epoch on the local training data, with a weighted loss to account for class imbalance. For the shared part of the model,  $f_{shared}$ , uniform weighting of the updates across clients is performed. Every experiment is conducted on the same 5 random seeds, leading to identical initialization and train-test splitting. The code, implemented in PyTorch [21], replicating the results on the public datasets, are available in a public code repository.<sup>1</sup> Further details on the hyperparameters, pre-processing and the model architecture can be found in Appendix A.3.

## 5 Results

### 5.1 Performance

From the perspective of a client choosing whether to participate in the federation, predictive performance is critical. In particular, attaining a collaborative model *worse* than a locally trained one virtually rules out participation. Therefore, two aspects are interesting: the average metric across all clients and the worst-performing client in the federation.

Holistically, iFedAvg shows competitive performance compared with a state-of-the-art benchmark in distributed personalized learning algorithm, APFL, and even outperforms both APFL, local and centralized training in several instances. As anticipated, in the case of heterogeneous or non-IID datasets, such as Ebola diagnosis and VSN, iFedAvg vastly outperforms FedAvg. Personalization appears important to adapt to these settings. Table 1 shows the average F1 score for each dataset and algorithm.

Analyzing the worst-performing client highlights the low tail of the performance distribution. In Table 2 we can observe that FedAvg performs especially poorly in this worst case. Intuitively, the client with a significantly shifted data distribution is *forced* to share the same global model, leading to inferior performance. This experiment highlights that iFedAvg is an effective method to personalize collaborative learning and is especially robust for the worst performing client in the federation. Additional tables and visualizations highlighting the full distribution of client performances, seed variation and additional metrics can be found in Appendix A.1.

---

<sup>1</sup>[github.com/davidroschewitz/ifedavg](https://github.com/davidroschewitz/ifedavg)

Table 1: Mean (across all clients) performance (F1 score)

| Dataset                    | Method         |       |        |       |              |
|----------------------------|----------------|-------|--------|-------|--------------|
|                            | iFedAvg (ours) | APFL  | FedAvg | Local | Centralized  |
| Ebola Prognosis            | 0.669          | 0.653 | 0.670  | 0.662 | <b>0.673</b> |
| Ebola Diagnosis            | <b>0.867</b>   | 0.828 | 0.773  | 0.844 | 0.861        |
| Vehicle Sensor Network     | 0.928          | 0.935 | 0.871  | 0.939 | <b>0.943</b> |
| Human Activity Recognition | <b>0.994</b>   | 0.992 | 0.967  | 0.993 | 0.988        |

Table 2: Worst-performing client in federation (F1 score)

| Dataset                    | Method         |       |        |              |             |
|----------------------------|----------------|-------|--------|--------------|-------------|
|                            | iFedAvg (ours) | APFL  | FedAvg | Local        | Centralized |
| Ebola Prognosis            | <b>0.581</b>   | 0.546 | 0.575  | 0.560        | 0.541       |
| Ebola Diagnosis            | <b>0.790</b>   | 0.662 | 0.452  | 0.654        | 0.733       |
| Vehicle Sensor Network     | 0.874          | 0.874 | 0.398  | <b>0.884</b> | 0.834       |
| Human Activity Recognition | 0.950          | 0.972 | 0.909  | <b>0.974</b> | 0.955       |

## 5.2 Interpretability of iFedAvg

One of the key improvements our method provides is that each locally personalized layer is directly interpretable. While the absolute value of each weight does not itself imply an exact relationship with the underlying data, the magnitude, direction and, most importantly, the comparison with other clients’ local weights, provides unparalleled insights into the interoperability of the datasets in the federation.

For the most valuable interpretable results, all clients are assumed to be willing to share the weights of their locally trained personalized layers,  $\mathbf{b}_{in}$  and  $\mathbf{w}_{in}$  for our experiments. In practice, this is a reasonable expectation, as raw data is not revealed and the insights might be critical for the clients.

In order to visually inspect the personalized shifts, we combine the weights into a heatmap for every feature and client to create a 2D matrix for each bias and weight. In the heatmap, we consider two types of shifts as *significant*. First, when for a single feature, a client’s local weight differs from the mean by more than 2 standard deviations (SD). Second, if the average SD of a feature in question differs by more than 2 SD from the average SD across *all features*. We show an example of such a heatmap on the VSN dataset in Figure 2, with the corresponding shift in the underlying data. Noteworthy is that our model detected the shift and indicated its direction without access to data of other clients, unlike the diagnosis histogram.

We further highlight the following examples of real-world shifts detected by our method:

- **Ebola Prognosis:** In Figure 3, we show the heatmap for patient prognosis. Every client significantly modifies the *CT Value* feature compared to other features. We can observe ETC Freetown’s distribution being more concentrated, and is therefore scaled by iFedAvg to increase compatibility with the shared model.
- **Ebola Diagnosis:** Similarly, we can observe effects for categorical differences. Figure 4 highlights that for ETC Foya, having the symptom of diarrhea is a clear indicator of positive Ebola infection (right in Figure). However this is confounded by systematic missingness for the diarrhea feature amongst EVD negative patients. iFedAvg corrects and identifies this difference in data collection without ever sharing any underlying data (left in Figure).

In summary, we have shown that iFedAvg is able to detect and correct various types of data shifts across clients in a collaborative learning setting that would ordinarily cause hidden bias and confounding. For highly non-IID datasets, we find large shift compensations (i.e., many significant values in the heatmap visualizations). This behavior is expected, but could be an indication of a non negligible false-positive rate. Given the objective of highlighting *potentially problematic* datasets and features, we believe this tradeoff is acceptable.

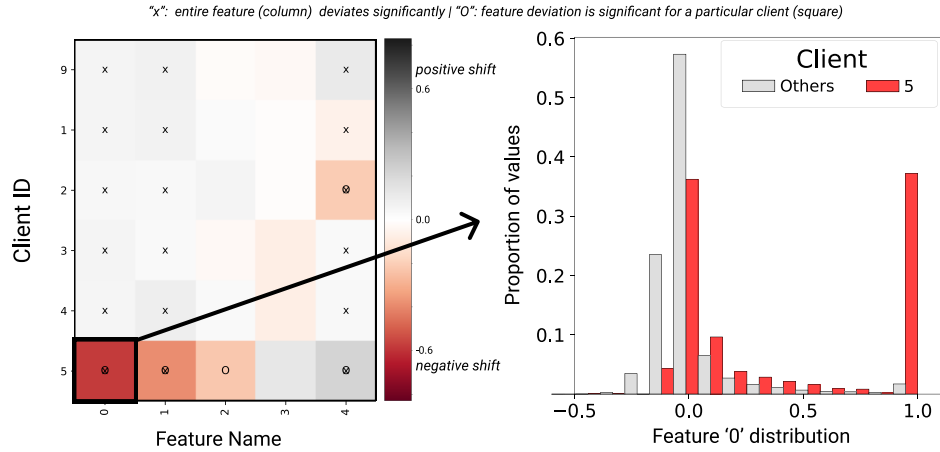


Figure 2: Input bias ( $w_{in}$ ) heatmap across clients for the VSN dataset. A negative shift is detected for sensor number 5 in feature '0' (left), which can be confirmed in the underlying distribution (right).

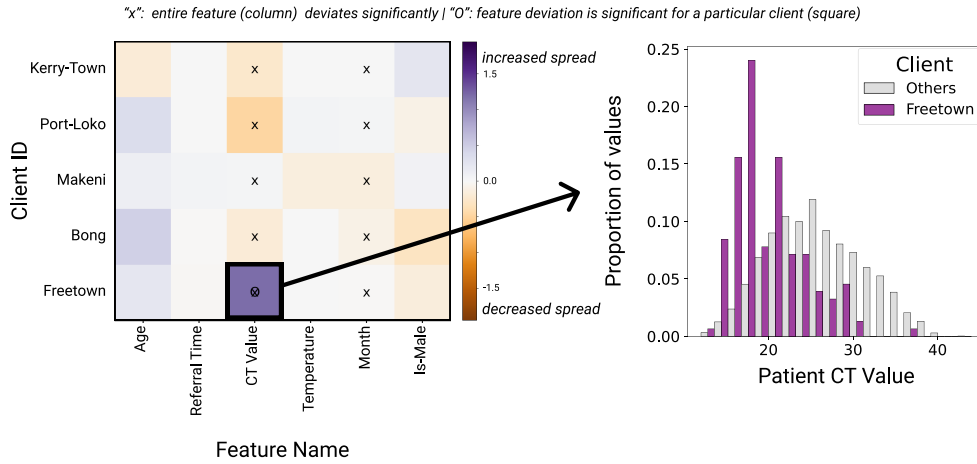


Figure 3: Heatmap of local input weights ( $w_{in}$ ) for Ebola patient prognosis (left) and underlying distribution of CT values for ETC 'Freetown' compared to all remaining clients (right).

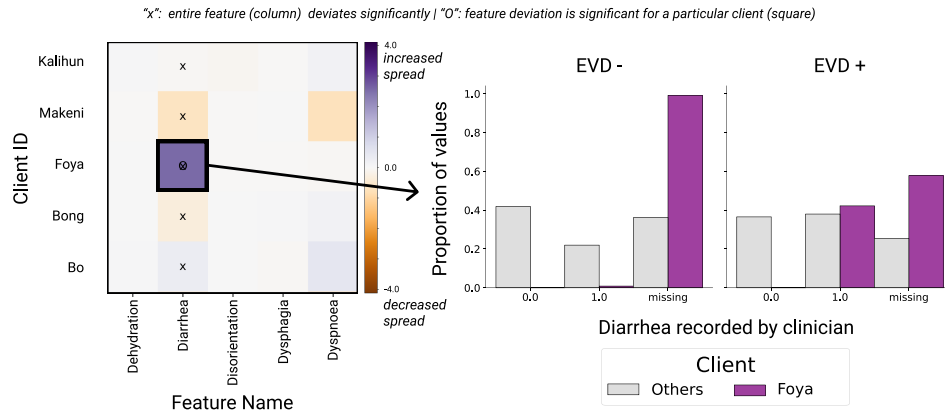


Figure 4: Local input weights heatmap ( $w_{in}$ ) for EVD diagnosis (left) and proportion of recorded diarrhea in the patient, split by final diagnosis for ETC 'Foya' compared to all remaining clients (right).

Interestingly, some of the personalized shifts detectable were indicative of data collection bias. Without a mechanism to flag such a shift, a personalized model trained by other methods would learn biased insights from data missing not-at-random (MNAR) and create poorly generalizable predictions. For example, let us recall the data collection of diarrhea presence for ETC Foya (Figure 4). A locally trained or strongly personalized federated model would achieve excellent predictive performance as patients with diarrhea appear to be correlated with Ebola infection. This pattern does not hold overall, and if falsely learned could lead to misdiagnosis without detection. In instances precisely like this, iFedAvg proves invaluable in detecting local *model biases*.

For various additional experimental results on the benchmark datasets as well as artificially introduced shifts and mutations confirming the efficacy of our method we refer to Appendix A.2.

## 6 Conclusion

We study interpretability and inter-client interoperability in federated learning and show how a feature-wise personalized learning approach addresses this challenge. Our framework, iFedAvg, proposes a simple extension to federated averaging that creates interpretable data-interoperability between clients by personalizing models. The learned weights ultimately reveal novel insights about the federated learning process as a whole.

On real-world datasets, iFedAvg is competitive with state-of-the-art personalized learning methods in terms of performance. The method vastly outperforms FedAvg or centralized learning on poorly performing clients with significant data shifts. More critically, significant feature-wise shifts in the underlying client datasets are correctly detected and compensated for. Not only does this provide targeted guidance to practitioners interpreting the results, but it can aid assessments of the overall compatibility of a client with the federation. While some shifts might not be harmful to the model, overcompensating for local data biases can be detrimental. iFedAvg visualizes these shifts, therefore generating the necessary transparency to best verify the reliability of model.

We leave as future research the extension of our approach to other types of data (such as images) or learning a feature alignment mapping. Furthermore, studying iFedAvg with a large number of clients and partial participation would allow more widespread adoption in the future.

In conclusion, iFedAvg offers novel insights into the federated learning process of tabular datasets. It leverages these insights not only for interpretability, but to build personalized models that are further adjusted for the usually hidden interoperability issues between clients in a federation. These unique extensions of the federated learning process come at a negligible computational overhead and thus iFedAvg is a promising approach for real world collaborative learning.

### Broader impact

This work is specifically designed to provide more guarantees of interoperability in federated learning and therefore incentivize collaboration, especially in fields with sensitive data and a high risk of collection bias. Equally we could disincentivize interoperable data collection by creating a shortcut that may undermine standardization efforts. Interoperability is a pillar of the FAIR Guiding Principles for ethical scientific data stewardship, and this critical issue was highlighted as a key barrier by a WHO-commissioned investigation into the massive failings of the centralized Ebola response [26]. Our work is specifically motivated by this use-case and is appropriately evaluated on a unique dataset collated from the largest number of distributed data collection sites during the notoriously poorly coordinated 2014-16 West African Ebola epidemic. By extrapolation, iFedAvg can be seen as a first step to facilitating interpretable interoperability in collaborative analyses. Basing the architecture on a distributed learning system, we also attempt to address the issue of local data ownership and data privacy compared to centralized approaches. While the insights shared do not reveal sample level data, the client-level aggregate could be considered sensitive and may be abused to discriminate against clients. In this instance, concealing the identity of the client could be considered as a mitigation strategy. Finally our simplified approach with low computational overhead makes this an accessible method for low-resource settings to build collaborative models whilst better securing patient privacy, intellectual property and statistical robustness to biases between collaborating datasets.

## Acknowledgements

The authors would like to acknowledge all the patients whose data was used in this study. This work was inspired by the challenges of sensitive data management in health emergencies and using the Ebola dataset provided critical validation and context. The data was provided by the Ebola Data Platform hosted by the Infectious Diseases Data Observatory (IDDO), and the data contributors, who had no role in the production of these research outputs. The contributors are: Alliance for International Medical Action (ALIMA), International Medical Corps (IMC), Institute of Tropical Medicine Antwerp (ITM), Médecins Sans Frontières (MSF), Oxford University and Save the Children (SCI). We also thank Aiyu Liu for his work in preparing the Ebola dataset for analysis and Mélanie Bernhardt for her critical and insightful thoughts and generous support.

## References

- [1] D. Anguita, A. Ghio, L. Oneto, X. Parra, and J. L. Reyes-Ortiz. A public domain dataset for human activity recognition using smartphones. In *Esann*, volume 3, page 3, 2013.
- [2] M. G. Arivazhagan, V. Aggarwal, A. K. Singh, and S. Choudhary. Federated learning with personalization layers. *arXiv preprint arXiv:1912.00818*, 2019.
- [3] A. Colubri, M. A. Hartley, M. Siakor, V. Wolfman, A. Felix, T. Sesay, J. G. Shaffer, R. F. Garry, D. S. Grant, A. C. Levine, and P. C. Sabeti. Machine-learning Prognostic Models from the 2014–16 Ebola Outbreak: Data-harmonization Challenges, Validation Strategies, and mHealth Applications. *EClinicalMedicine*, 11:54–64, may 2019. ISSN 25895370. doi: 10.1016/j.eclinm.2019.06.003.
- [4] L. Corinzia, A. Beuret, and J. M. Buhmann. Variational federated multi-task learning. *arXiv preprint arXiv:1906.06268*, 2019.
- [5] Y. Deng, M. M. Kamani, and M. Mahdavi. Adaptive personalized federated learning. *arXiv preprint arXiv:2003.13461*, 2020.
- [6] M. F. Duarte and Y. H. Hu. Vehicle classification in distributed sensor networks. *Journal of Parallel and Distributed Computing*, 64(7):826–838, 2004.
- [7] A. Fallah, A. Mokhtari, and A. Ozdaglar. Personalized federated learning: A meta-learning approach. *arXiv preprint arXiv:2002.07948*, 2020.
- [8] Georgetown University Medical Center. Data Sharing during the West Africa Ebola Public Health Emergency: Case Study Report. Technical report, Georgetown University Medical Center, 2018.
- [9] M. A. Hartley, A. Young, A. M. Tran, H. H. Okoni-Williams, M. Suma, B. Mancuso, A. Al-Dikhari, and M. Faouzi. Predicting Ebola Severity: A Clinical Prioritization Score for Ebola Virus Disease. *PLoS Neglected Tropical Diseases*, 11(2):e0005265, feb 2017. ISSN 19352735. doi: 10.1371/journal.pntd.0005265.
- [10] M. A. Hartley, A. Young, A. M. Tran, H. H. Okoni-Williams, M. Suma, B. Mancuso, A. Al-Dikhari, and M. Faouzi. Predicting Ebola infection: A malaria-sensitive triage score for Ebola virus disease. *PLoS Neglected Tropical Diseases*, 11(2):e0005356, feb 2017. ISSN 19352735. doi: 10.1371/journal.pntd.0005356.
- [11] T.-M. H. Hsu, H. Qi, and M. Brown. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- [12] A. Imakura, H. Inaba, Y. Okada, and T. Sakurai. Interpretable collaborative data analysis on distributed data. *Expert Systems with Applications*, 177:114891, 2021.
- [13] Infectious Disease Data Observatory (IDDO). Ebola Data Platform, 2021. URL <https://www.iddo.org/research-themes/ebola>.

- [14] V. Jain, A. Charlett, and C. S. Brown. Meta-analysis of predictive symptoms for ebola virus disease. *PLoS Neglected Tropical Diseases*, 14(10):1–15, oct 2020. ISSN 19352735. doi: 10.1371/journal.pntd.0008799.
- [15] Y. Jiang, J. Konečný, K. Rush, and S. Kannan. Improving federated learning personalization via model agnostic meta learning. *arXiv preprint arXiv:1909.12488*, 2019.
- [16] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- [17] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.
- [18] V. Kulkarni, M. Kulkarni, and A. Pant. Survey of personalization techniques for federated learning. In *2020 Fourth World Conference on Smart Trends in Systems, Security and Sustainability (WorldS4)*, pages 794–797. IEEE, 2020.
- [19] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith. Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.
- [20] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pages 1273–1282. PMLR, 2017.
- [21] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library, 2019.
- [22] S. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan. Adaptive federated optimization. *arXiv preprint arXiv:2003.00295*, 2020.
- [23] M. T. Ribeiro, S. Singh, and C. Guestrin. Model-agnostic interpretability of machine learning. *arXiv preprint arXiv:1606.05386*, 2016.
- [24] V. Smith, C.-K. Chiang, M. Sanjabi, and A. Talwalkar. Federated multi-task learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 4427–4437, 2017.
- [25] G. Wang. Interpret federated learning with shapley values. *arXiv preprint arXiv:1905.04519*, 2019.
- [26] M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L. B. da Silva Santos, P. E. Bourne, et al. Comment: The fair guiding principles for scientific data management and stewardship. *Scientific data*, 3:160018, 2016.
- [27] T. Yu, E. Bagdasaryan, and V. Shmatikov. Salvaging federated learning by local adaptation. *arXiv preprint arXiv:2002.04758*, 2020.
- [28] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.
- [29] F. Zheng, K. Li, J. Tian, X. Xiang, et al. A vertical federated learning method for interpretable scorecard and its application in credit scoring. *arXiv preprint arXiv:2009.06218*, 2020.

## A Appendix

### A.1 Supplementary performance results

In addition to the main results in the paper, we provide the full distribution of performance metrics across clients as well as balanced accuracy and ROC AUC scores for each method. In the client performance distribution plots, the red error bars show the standard deviation (SD) of the median score of each random seed. We notice no particular pattern or single method with a distinctive behavior with regards to seed. The following two sections evaluate this experiment by F1 score and AUROC and .

#### A.1.1 F1 Score

Tables 1 and 2 (in the main text) show mean and worst-client performance on all datasets. Here, we show the distribution of client performances in violin plots for each dataset (EVD Prognosis: Figure 5, EVD Diagnosis: Figure 6, VSN: Figure 7, HAR: Figure 8).

It is particularly apparent in these visualizations that *iFedAvg* is robust to *poorly* performing clients. For every dataset our method outperforms FedAvg, and in most instances even outperforms APFL, Local or Centralized training. Furthermore, for the HAR dataset, while *iFedAvg* has a relatively low performing client, the overall distribution is skewed towards 1.0, indicating good overall performance. For this dataset, the seed-SD is also noticeably lower for our method compared to the benchmarks.

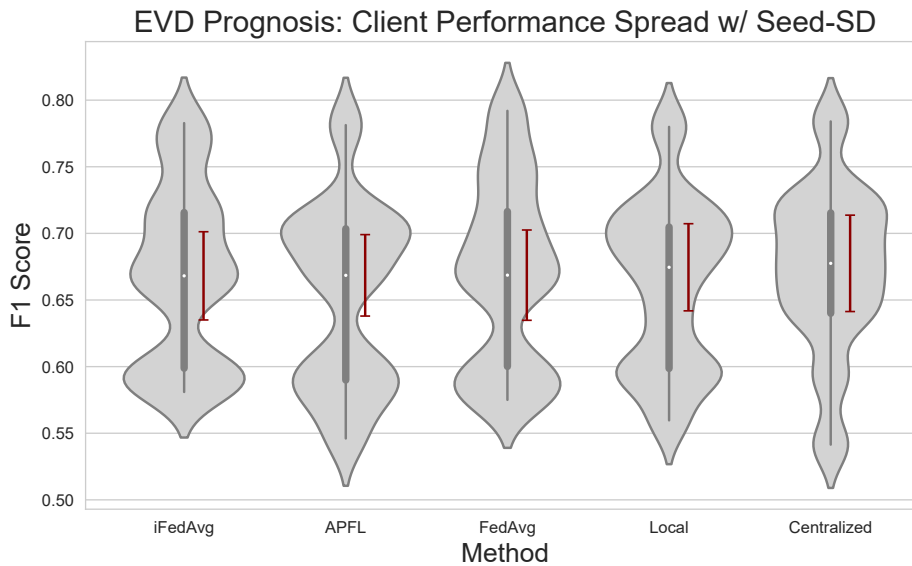


Figure 5: Distribution of client performances (F1 score) for the Ebola Prognosis dataset. Red error bars show the standard deviation (SD) of the median score of each random seed.





Figure 6: Distribution of client performances (F1 score) for the Ebola Diagnosis dataset. Red error bars show the standard deviation (SD) of the median score of each random seed.

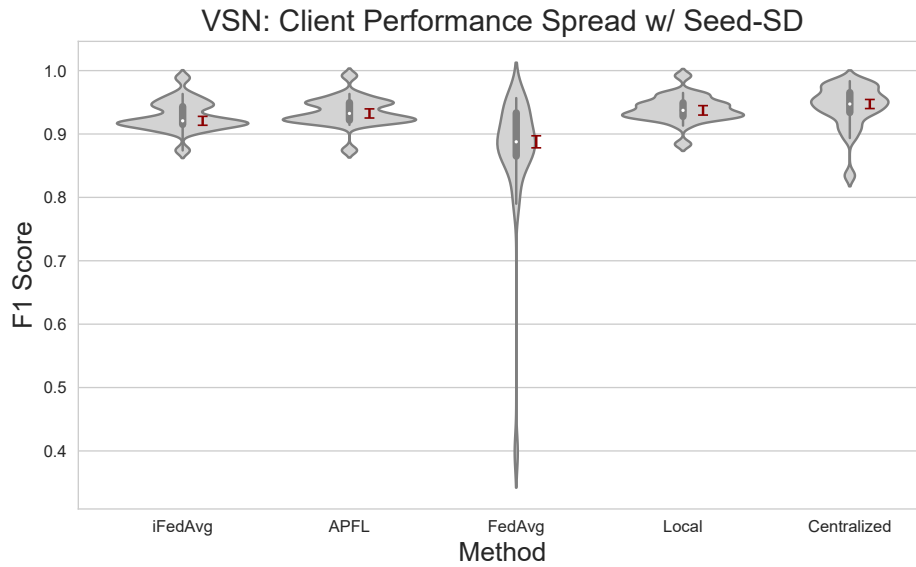


Figure 7: Distribution of client performances (F1 score) for the Vehicle Sensor Network dataset. Red error bars show the standard deviation (SD) of the median score of each random seed.

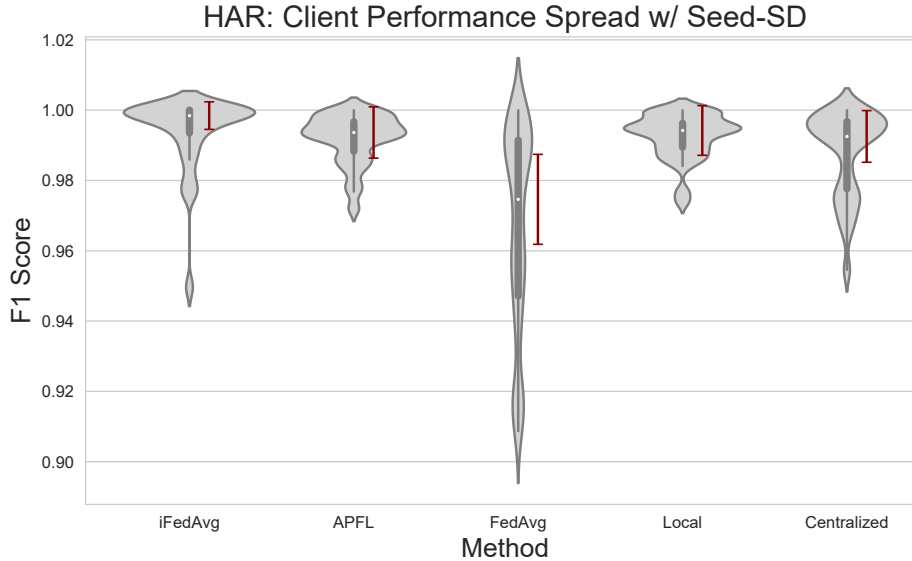


Figure 8: Distribution of client performances (F1 score) for the Human Activity Recognition dataset. Red error bars show the standard deviation (SD) of the median score of each random seed.

Table 3: Mean (across all clients) performance (ROC AUC)

| Dataset                    | Method         |       |              |       |              |
|----------------------------|----------------|-------|--------------|-------|--------------|
|                            | iFedAvg (ours) | APFL  | FedAvg       | Local | Centralized  |
| Ebola Prognosis            | 0.725          | 0.704 | <b>0.726</b> | 0.708 | 0.74         |
| Ebola Diagnosis            | <b>0.909</b>   | 0.860 | 0.879        | 0.870 | 0.901        |
| Vehicle Sensor Network     | 0.975          | 0.978 | 0.921        | 0.982 | <b>0.986</b> |
| Human Activity Recognition | 1.00           | 1.00  | 0.999        | 1.00  | 1.00         |

### A.1.2 ROC AUC

An additional metric which can be used to measure the performance of a classification model is the area under the curve of the receiver operating characteristic. We evaluate our method in the same fashion as previously; analyzing both the average performance across clients as well as the worst performing client in the federation. These results are displayed in Table 3 and 4.

These graphs corroborate the findings in the main text, and highlight that the strong performance of iFedAvg is independent of the chosen metric. Interestingly, the gap between our method and vanilla federated averaging shrinks marginally, which can be explained by the sensitivity of F1 score to individual incorrect samples.

Similarly to the F1 score metric, we present all distributions of ROC AUC performance as violin plots (EVD Prognosis: Figure 9, EVD Diagnosis: Figure10, VSN: Figure 11, HAR: Figure 12).

Table 4: Worst-performing client in federation (ROC AUC)

| Dataset                    | Method         |       |        |              |              |
|----------------------------|----------------|-------|--------|--------------|--------------|
|                            | iFedAvg (ours) | APFL  | FedAvg | Local        | Centralized  |
| Ebola Prognosis            | <b>0.628</b>   | 0.573 | 0.608  | 0.567        | 0.581        |
| Ebola Diagnosis            | <b>0.799</b>   | 0.694 | 0.770  | 0.689        | 0.740        |
| Vehicle Sensor Network     | 0.930          | 0.936 | 0.124  | <b>0.950</b> | 0.935        |
| Human Activity Recognition | 0.996          | 0.997 | 0.991  | 0.997        | <b>0.998</b> |

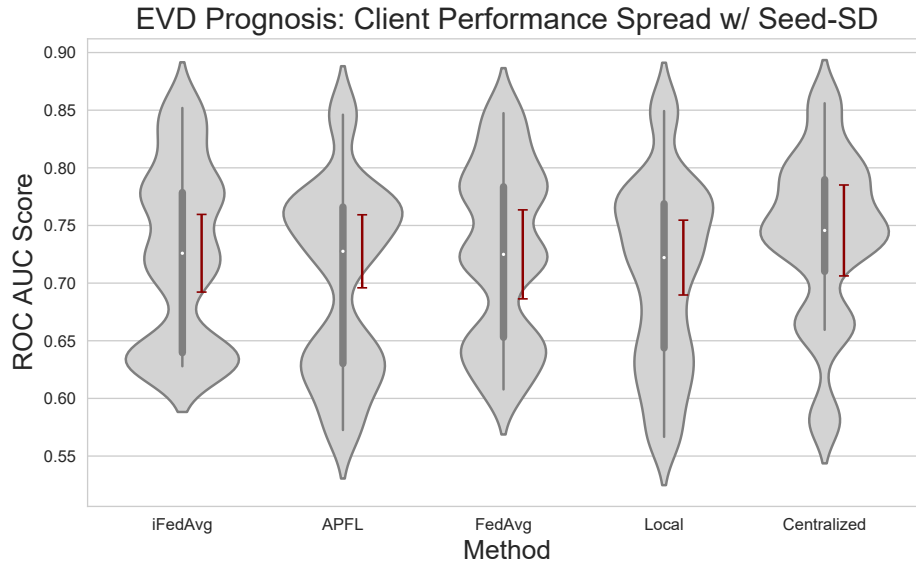


Figure 9: Distribution of client performances (ROC AUC) for the Ebola Prognosis dataset.

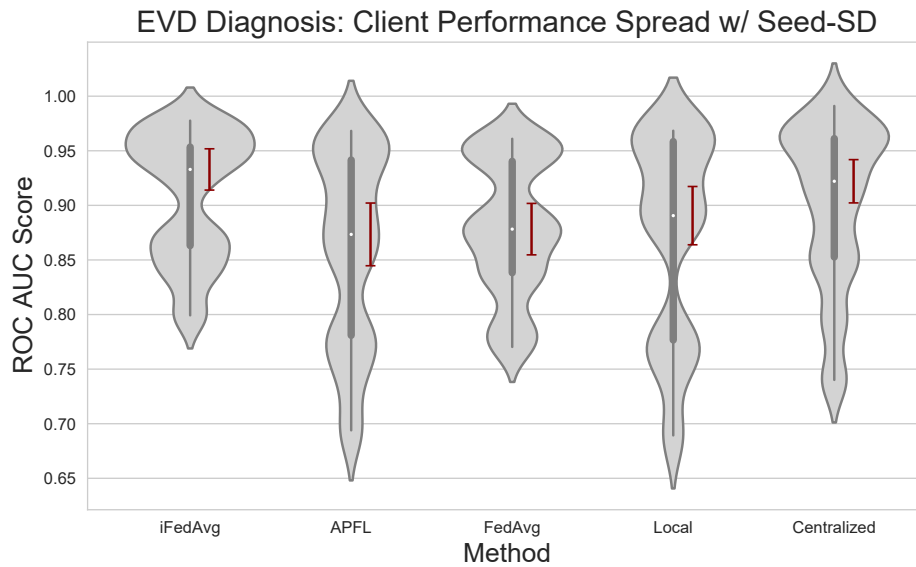


Figure 10: Distribution of client performances (ROC AUC) for the Ebola Diagnosis dataset.

One instance that stands out is the relatively poor performance of APFL for EVD Prognosis and Diagnosis, as the personalized method does not manage to outperform vanilla federated averaging. We suspect that APFL is overfitting due to the strongly heterogeneous nature of the datasets, which can be observed in its similarity to Local training performance.

In conclusion, iFedAvg shows impressive performance in various settings and measured by different metrics. Especially on datasets which benefit from personalization due to one or multiple outlier clients, our method is particularly effective compared with vanilla federated averaging.

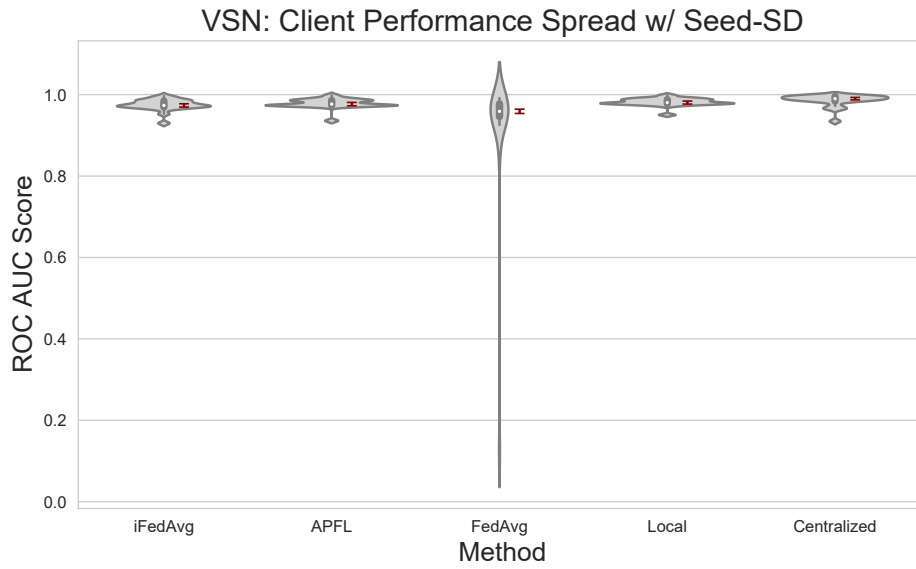


Figure 11: Distribution of client performances (ROC AUC) for the Vehicle Sensor Network dataset.

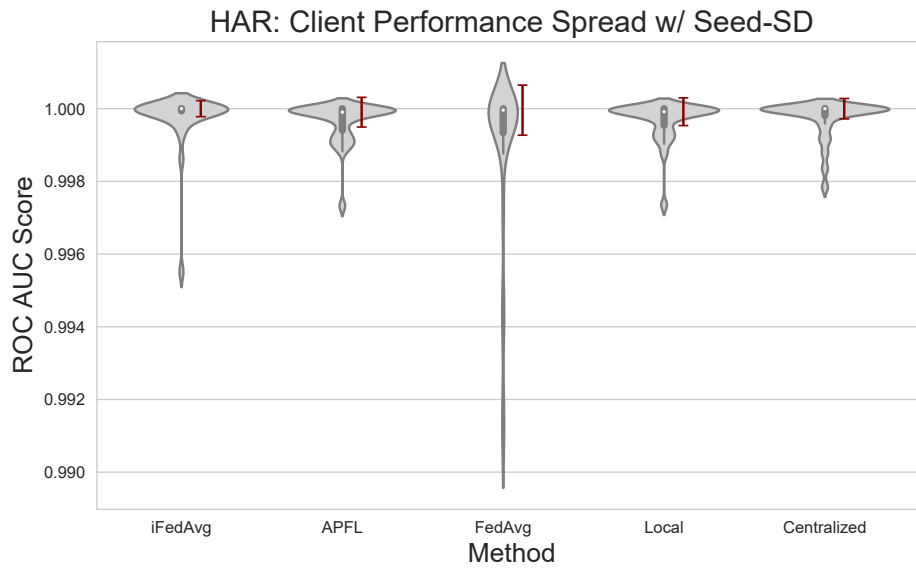


Figure 12: Distribution of client performances (ROC AUC) for the Human Activity Recognition dataset.

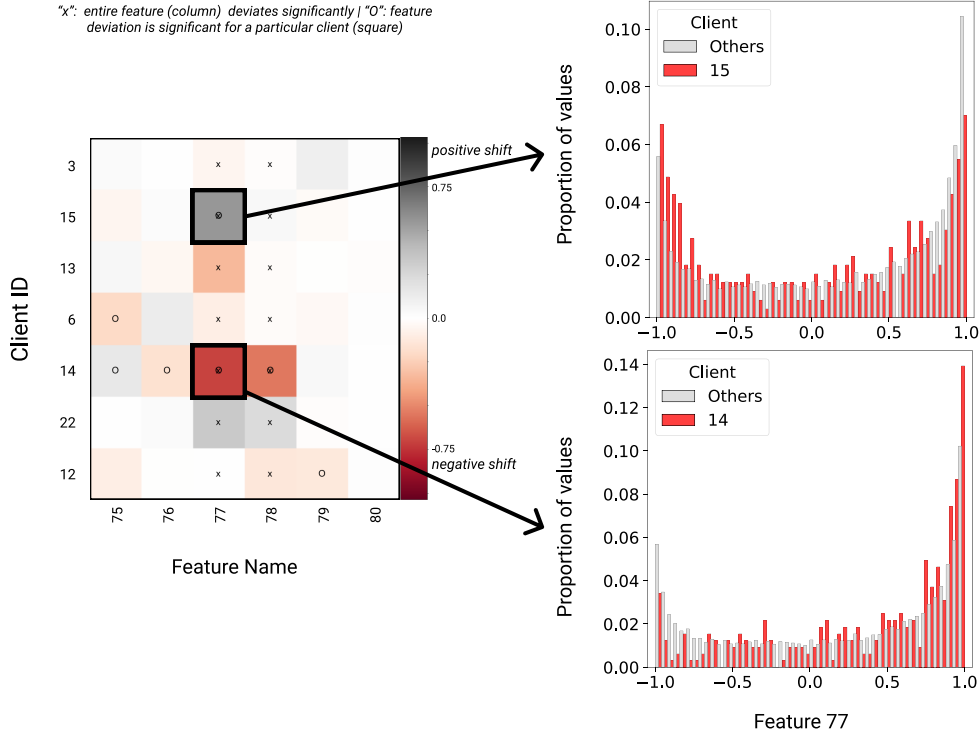


Figure 13: Heatmap of the biases ( $b_{in}$ ) for the Human Activity Recognition dataset (left) and histograms of feature 77 for clients 15 and 14 (right).

## A.2 Supplementary interpretability results

In this section, we show additional interpretability results of iFedAvg, with the objective of highlighting the different types of shifts detected for various datasets. The list below serves as a reference to the corresponding figures.

- Human Activity Recognition - underlying shifts:** In order to determine the directionality of detected shifts, we highlight two examples that are shown for iFedAvg. For feature 77, we notice that all shifts are significant ('x' in the column), but two local weights for clients 14 and 15 are additionally significant ('O'). We can see that the positive and negative shifts in bias ( $b_{in}$ ) in grey and red, respectively, are noticeable in the underlying data. The histograms clearly show that the values of client 15 are skewed towards  $-1.0$  whereas for client 14 they are skewed towards  $1.0$ . This example demonstrates that the learned local layers can indicate directionality correctly. Figure 13 shows the bias heatmap as well as the histograms of both clients of interest.
- Human Activity Recognition - feature spread:** Highlighting another example on the HAR dataset, our method detects that for feature 50, client 16, a smaller weight is being applied. Investigating the underlying feature shows that, for this client, larger values are not being observed. While personalizing the local layers, iFedAvg therefore is reducing the value of this feature. This could be indicative of *feature deactivation* or simply a compensation for the downstream effect of this feature in the *shared* part of the model. The histogram of the feature and heatmap of  $w_{in}$  can be seen in Figure 14.
- Vehicle Sensor Network - Target-specific Stretch:** A particularly interesting example of shifts detected by iFedAvg is one that is dependent on the target. For the VSN dataset, client 5 seems to have different values for feature 50 for the positive class, which is shown as significant in the heatmap. Client 19, albeit not significantly, has a slightly above average

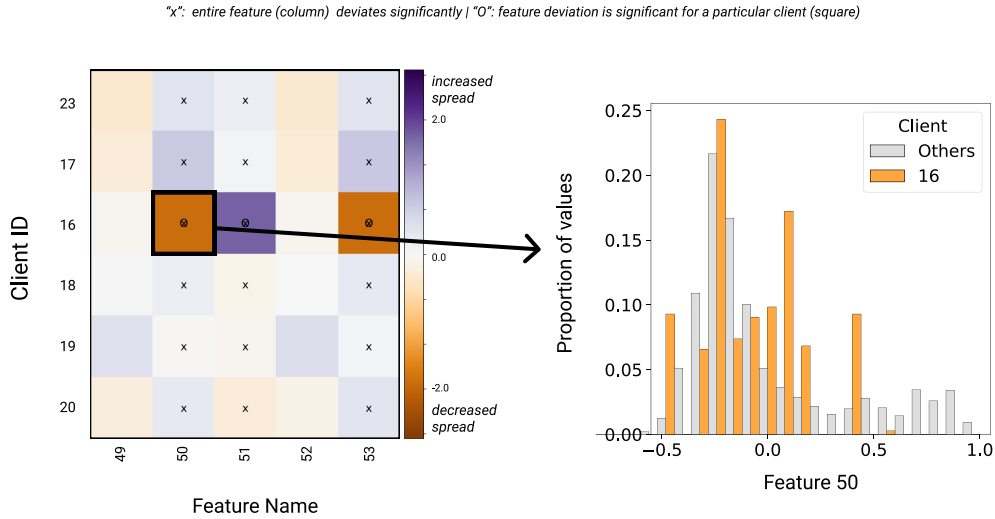


Figure 14: Heatmap of the weights ( $w_{in}$ ) for the Human Activity Recognition dataset (left) and histograms of feature 50 for client 16 (right).

weight, and indeed for the negative target class has a differing underlying distribution. While the directionality of the weights in this instance do not directly indicate a shift for a particular target class, our method correctly identifies *differences* in the underlying feature distribution. The heatmap and histograms of feature 50 for both clients and both target classes can be seen in Figure 15.

- **Ebola Diagnosis - small underlying shifts:** Some of the discussed examples have been rather substantial and for features where many clients modify the bias or weight (marked by 'x' in the columns). Detecting and correcting, small shifts can be critical for a client if they are the only ones performing such a personalized compensation. For the Kalihun ETC, it appears as if the referral times, the time taken until a patient actually visits the treatment center, are slightly larger and  $iFedAvg$  is compensating with a small but significant negative bias. For comparison, the histogram of an arbitrary ETC, Foya, is shown highlighting the difference in the underlying data. The heatmap and histograms are shown in Figure 16
- **Ebola Diagnosis - data collection differences:** As *bias* in the local data could have catastrophic consequences, we highlight another example. For ETC Foya, whether a patient has a malaria co-infection only appears to be recorded for EVD-positive cases. Here, the ETC only records a malaria test in confirmed cases. For comparison, another ETC, Kalihun, is shown, which does not record malaria infection at all. Without explicitly detecting this effect, the personalized model might overfit in practice, with detrimental consequences. The heatmap and both histograms are shown in Figure 17.

The previous examples show that  $iFedAvg$  is able to detect and compensate for various types of discrepancies in the underlying local datasets. Some, however, might not be desirable, and therefore being able to identify, at a client and feature level, problematic values could help reduce model bias.

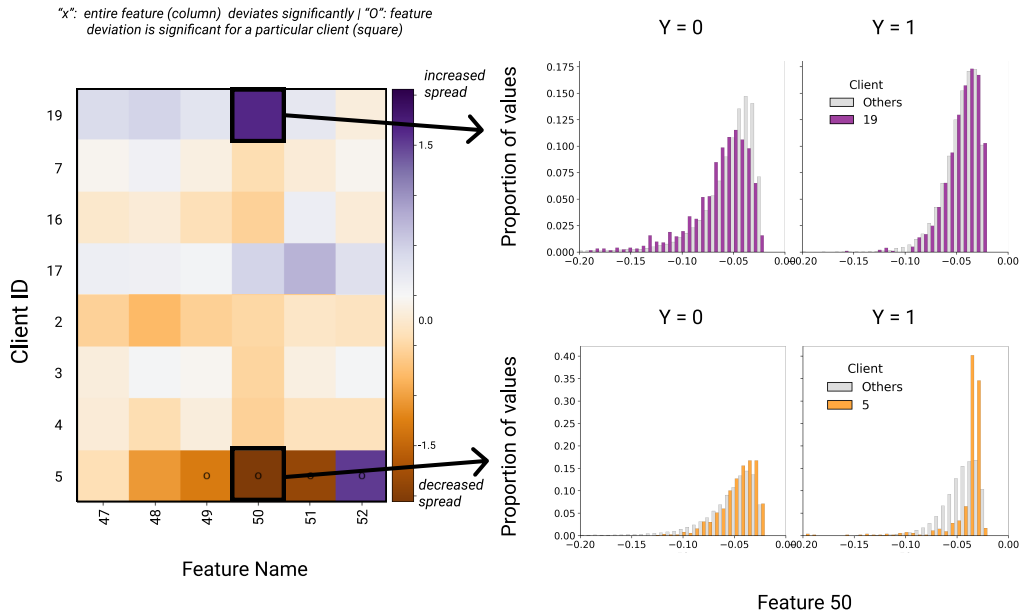


Figure 15: Heatmap of the weights ( $w_{in}$ ) for the Vehicle Sensor Network dataset (left) and histograms of feature 50 for clients 19 and 5 split according to the target class (right).

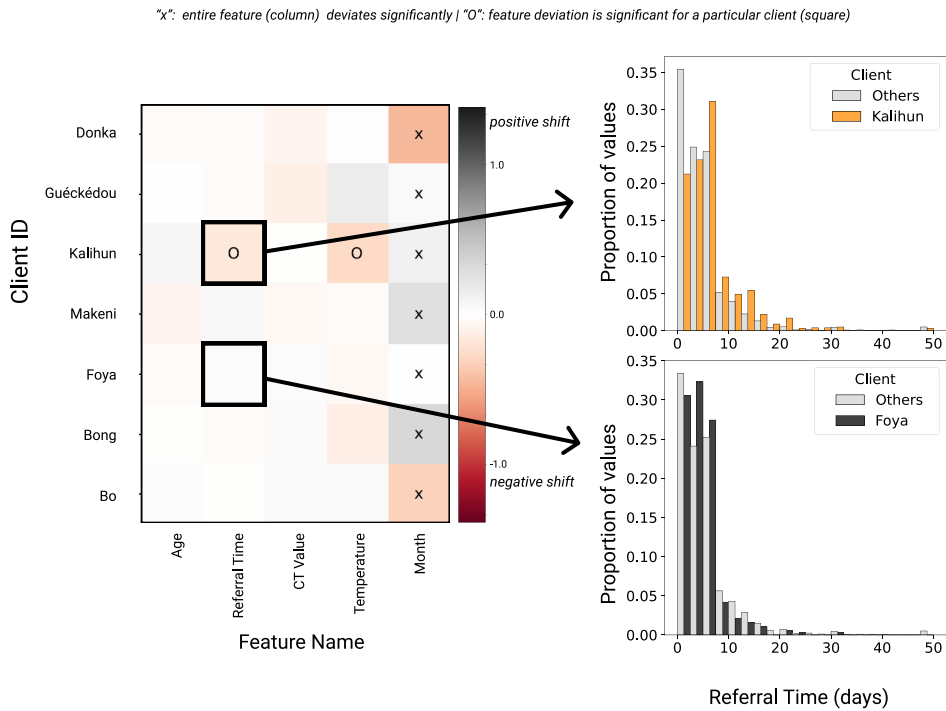


Figure 16: Heatmap of the biases ( $b_{in}$ ) for the Ebola Diagnosis dataset (left) and histograms of the referral time feature for ETCs Kalihun and Foya (right).

\*X\*: entire feature (column) deviates significantly | \*O\*: feature deviation is significant for a particular client (square)

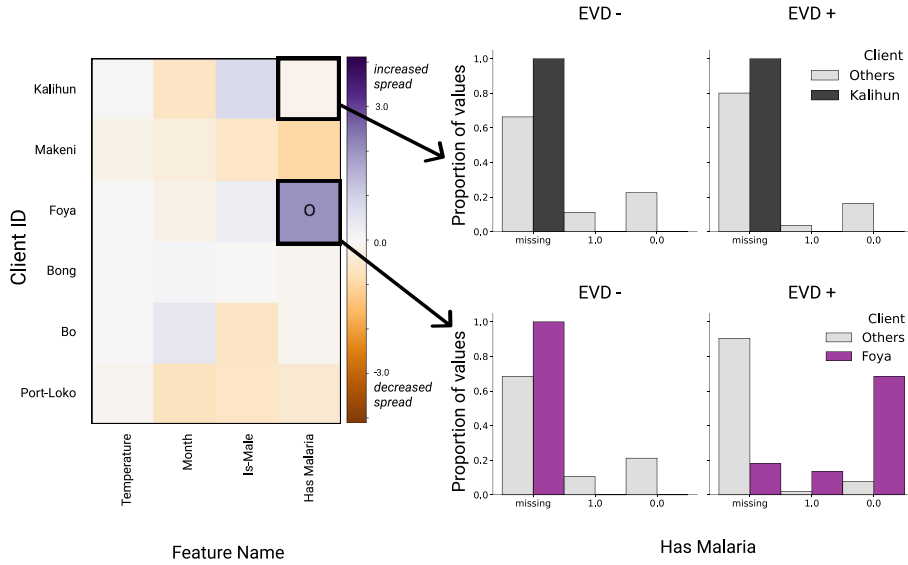


Figure 17: Heatmap of the weights ( $w_{in}$ ) for the Ebola Diagnosis dataset (left) and histograms of the Malaria infection feature for ETCs Kalihun and Foya, split according to the target class (EVD- and EVD+) (right).

Table 5: Architecture of the MLP model used

| Type            | Size (in, out) | Activation | Dropout |
|-----------------|----------------|------------|---------|
| $f_{in}$        | $(D, D)$       | -          | 0.2     |
| Fully-connected | $(D, 128)$     | TanH       | 0.2     |
| Fully-connected | $(128, 64)$    | TanH       | 0.2     |
| Fully-connected | $(64, K)$      | -          | -       |
| $f_{out}$       | $(K, K)$       | -          | -       |

### A.3 Hyperparameters and experimental setup

In order to create the most comparable experiments, an identical network architecture was used for all experiments. This ensures that each method has the same number of parameters available in the base model. APFL, of course, creates multiple copies of this model. We show the entire MLP architecture in Table 5, with  $D$  the number of features and  $K$  the output dimension.  $K$  is chosen to be the number of classes in our experiments, after which log-softmax is applied in conjunction with negative-log-likelihood loss. The class weights, used to weight the loss function, are computed as the inverse of class prevalence, scaled to sum to  $K$ .

Each experiment was conducted on the following seeds, which dictate the local train and holdout set splitting, network initialization and batch shuffling. 2934384, 10231938, 8273, 2019231, 62739. The learning rate of SGD was set to 0.002 for all experiments and datasets, as this performed best across the board. The learning rate was decayed using a step function 50 times, with a step of 0.9 for the entirety of the 1000 rounds. Client-side momentum with a value of 0.5 was enabled for all methods except APFL as the authors do not discuss it. For APFL the best  $\alpha$  was empirically found to be 0.5. The F1 score was computed in a weighted fashion, ROC AUC with one-vs-one treatment for multi-class targets.

For both Ebola datasets, each client locally standardized the numerical features. For the VSN and HAR datasets, the original standardization of the benchmark datasets was retained. While there is no significant difference in the results, both modes are supported by iFedAvg and implemented in the opensourced code.



Table 6: Performance difference with  $f_{\text{out}}$  enabled (F1 score)

| Dataset                    | $\Delta$ Average | $\Delta$ Worst-performing client |
|----------------------------|------------------|----------------------------------|
| Ebola Prognosis            | +0.002           | +0.008                           |
| Ebola Diagnosis            | +0.005           | +0.008                           |
| Vehicle Sensor Network     | +0.002           | -                                |
| Human Activity Recognition | -                | -0.013                           |

#### A.4 Target shift layer

Intuitively, the layer  $f_{\text{out}}$  acts as a personalized compensation of any differences in how the target differs for each client. For instance, one would hope to detect varying class imbalance and a less clear class distinction in this layer. Interpreting this layer is slightly less intuitive, as each value corresponds to a logit, not a real feature. Nonetheless we present results in this section of iFedAvg with the training of local  $f_{\text{out}}$  enabled.

First, we analyze the performance with the target layer enabled. As can be seen in Table 6, for most datasets there is a marginal performance gain. This difference is not significant enough to warrant this layer as necessary, but also highlights that it is not detrimental.

We structure our analysis into the following three sections:

- **Overall results:** We show the heatmaps of  $\mathbf{b}_{\text{out}}$  and  $\mathbf{w}_{\text{out}}$  for EVD Diagnosis in Figure 18. Each column now no longer represents input features, but the logit of each target class (0 being negative, 1 being positive diagnosis). In addition to the setting with both the bias and weight being trained locally, we also explore enabling the personalized training on each independently. It is visually apparent, that there are many smaller deviations, but few stand out except the weight of class 1 and ETC Foya. These shifts are not directly correlated with easily visible characteristics such as positivity rate. Therefore, the learned shifts are more complex to diagnose, and should be used as an indicator of potentially other issues. This is an interesting future research avenue.
- **Detecting mis-labeled targets:** While not as common as poorly standardized features, it is possible that due to a simple translation error, the label for a single client is flipped. With methods other than iFedAvg, this would result in terrible performance, and no interpretable output to help detect the issue. With an enabled target layer,  $f_{\text{out}}$ , this becomes quite trivial. Figure 19 shows the heatmap for EVD Diagnosis, where for ETC Kalihun, the label was artificially swapped. An interesting alternative definition of  $\mathbf{w}_{\text{out}}$  would be as a single scalar value, instead of a vector  $\in \mathbb{R}^K$ . This only allows each client a single multiplicative scaling of the last outputs. We show this setting in the right hand side of Figure 19. As can be observed in both constellations, a large negative weight is learned, indicating that for one class, the label is inverted. This is particularly apparent in the scalar setting. With this information, a client in the federation would have a strong indication that an interoperability issue does not originate from the features, but from the target itself.
- **Consistent feature weights:** An undesirable effect of enabling the personalized target layer learning would be an impact on the feature-wise interpretable results. We highlight how this effect is marginal in Figure 20. Therefore, the target layer can be enabled for most use-cases, especially when a target-swap could occur.

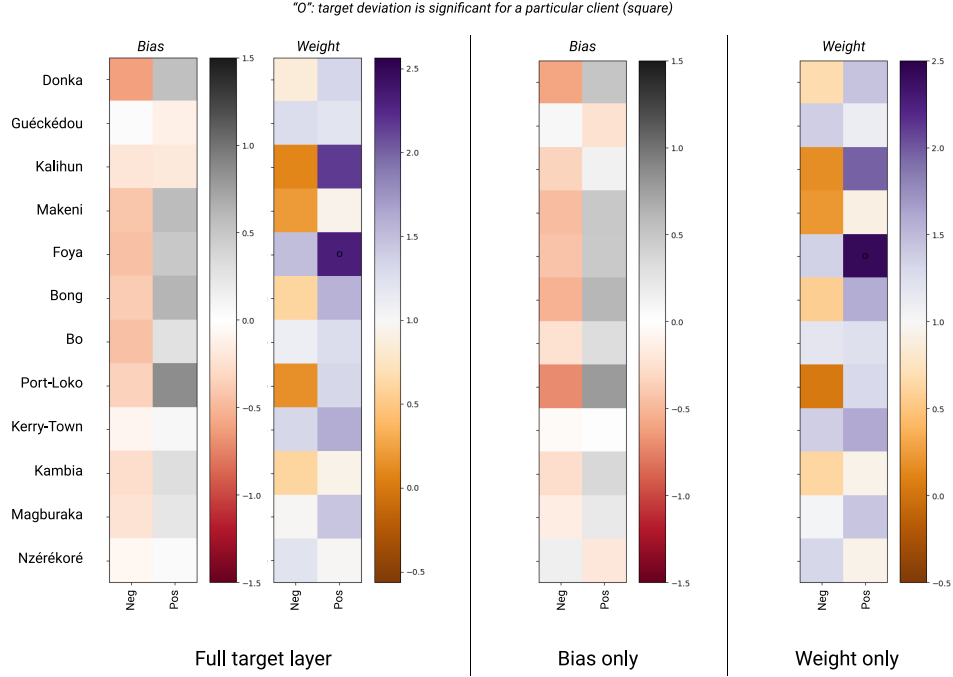


Figure 18: Heatmap of biases and weights ( $\mathbf{b}_{out}$  and  $\mathbf{w}_{out}$ ) of the target layer for the Ebola Diagnosis dataset. Each column represents the positive and negative class. Three configurations are shown: 1) both bias and weight trained locally (left), 2) only bias trained locally (middle) and only weight trained locally (right).

### A.5 Ebola dataset statistics

In Table 7 and 8 we show the number of samples at each treatment center. For prognosis, only patients where the outcome is known and a patient was confirmed EVD positive are considered. For diagnosis, only patients where an EVD test was performed and the minority class has at least 2% of samples are included. This leads to the fact that not the same ETCs are represented for both tasks. The reason being that some ETCs did not monitor mortality, or others only treated EVD+ cases.

Table 7: Ebola Prognosis dataset summary statistics

| ETC            | Number of samples | Mortality rate |
|----------------|-------------------|----------------|
| Guéckédou      | 1366              | 66.98%         |
| Monrovia       | 1154              | 56.85%         |
| Kalihun        | 852               | 44.37%         |
| Donka (EJPDEJ) | 748               | 49.87%         |
| Foya           | 450               | 66.00%         |
| Bo             | 440               | 38.63%         |
| Donka (EFFVXT) | 418               | 37.80%         |
| Kerry-Town     | 263               | 42.59%         |
| Port-Loko      | 181               | 65.75%         |
| Makeni         | 176               | 56.82%         |
| Bong           | 168               | 50.00%         |
| Freetown       | 166               | 50.00%         |

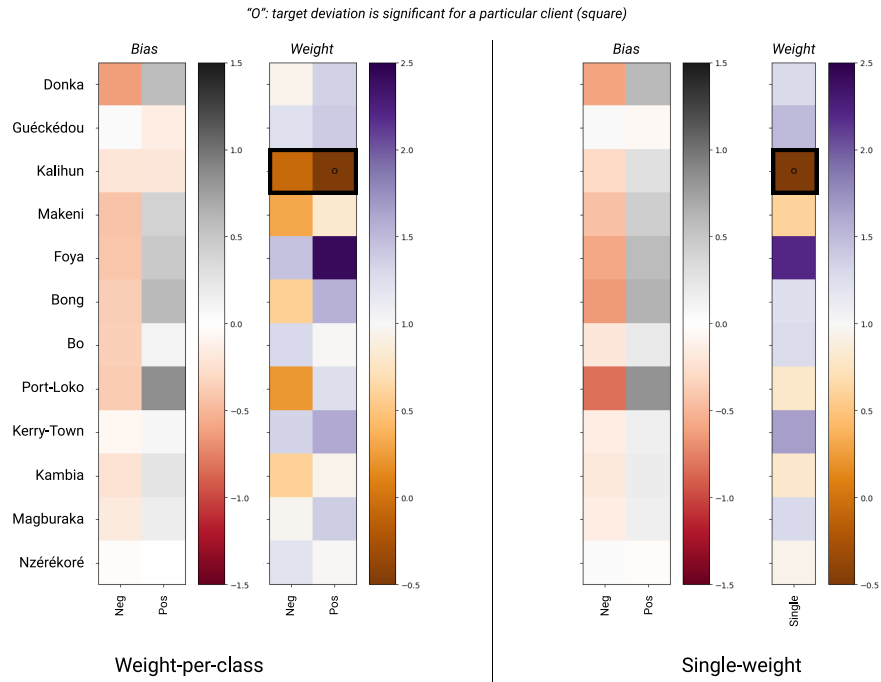


Figure 19: Heatmap of biases and weights ( $\mathbf{b}_{\text{out}}$  and  $\mathbf{w}_{\text{out}}$ ) of the target layer for the Ebola Diagnosis dataset with an artificially introduced target flip for ETC Kalihun (marked). Two configurations are shown: 1)  $\mathbf{w}_{\text{out}}$  is a vector, with a value for each target class (left) and  $\mathbf{w}_{\text{out}}$  is a scalar with a single value (right).

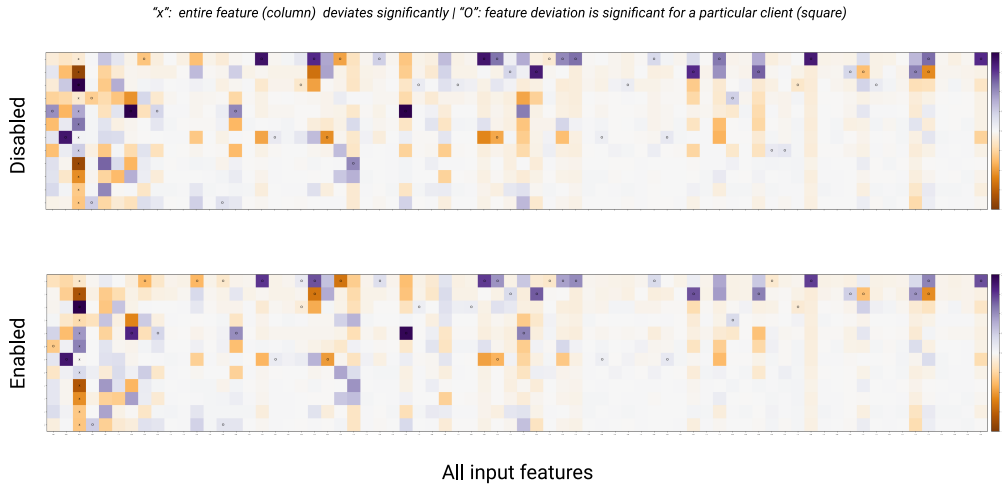


Figure 20: Heatmaps of the feature weights ( $\mathbf{w}_{\text{in}}$ ) for the Ebola Diagnosis dataset with local training of the target layer (bias and weight) disabled (top) and enabled (bottom)

Table 8: Ebola Diagnosis dataset summary statistics

| ETC        | Number of samples | Positivity rate |
|------------|-------------------|-----------------|
| Donka      | 1975              | 37.87%          |
| Guéckédou  | 1517              | 90.05%          |
| Kalihun    | 1173              | 72.63%          |
| Makeni     | 848               | 20.75%          |
| Foya       | 564               | 79.79%          |
| Bong       | 529               | 31.76%          |
| Bo         | 519               | 84.78%          |
| Port-Loko  | 477               | 37.95%          |
| Kerry-Town | 275               | 95.64%          |
| Kambia     | 217               | 21.66%          |
| Magburaka  | 155               | 29.03%          |
| Nzérékoré  | 137               | 57.66%          |