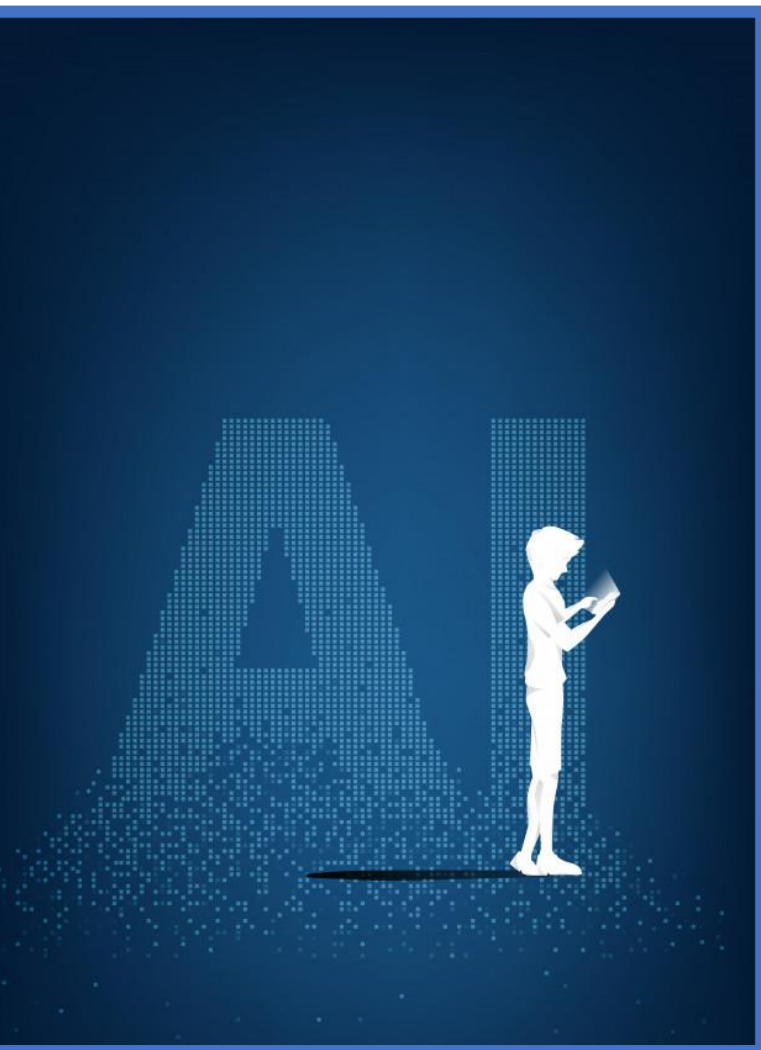


Practical Byzantine-resilient, yet decentralized federated learning

A thesis written by Joost Verbraeken examining the state-of-the-art and proposing Pro-Bristle, a new technique to improve byzantine-resilience in asynchronous non-i.i.d. settings



Abstract

Introduction

(glossary)

Why machine learning?

Statistics is a branch of mathematics concerned with explain and gathering insights from historical data, for example to understand consumer preferences, determine the core components of human personalities, or to illustrate how economic policies affect the society. However, statistics has its limitations: it can be notoriously hard to create a highly accurate model, the statistical techniques available for making predictions about the future are relatively limited, and the use-cases are limited to clearly specified domains (in contrast to fuzzy domains such as the generation of completely new songs based on previous songs). Machine learning gained traction over the last few decades thanks to the discovery of several novel and powerful techniques, such as Support Vector Machines and Random Forests. These techniques are used for a wide variety of tasks such as multimedia recommendation, handwriting recognition, speech-to-text systems, digital translators, etc., but these methods depend on carefully extracted features and often yield sub-optimal accuracy. In the last decade, neural networks became popular thanks to their versatility, ability to learn to extract proper features themselves, and highly effective predictions. Neural networks consist of a series of layers, each with a number of artificial neurons. Each neuron is linked to a number of other neurons in the previous/next layer by a connection associated with a certain weight. When a neuron is activated, it checks if the combined activation it gets from all nodes in the last layer exceeds a certain threshold (i.e. its bias) and then propagates this activation to the nodes in the next layer to which it is connected. These weights and biases are constantly adjusted in the neural network through a process called back-propagation so that the network actually starts to “learn”.

Why distributed learning?

Training a neural network on a single machine is possible when the amount of data is relatively limited. However, for more complex applications (such as self-driving cars, image recognition, or music generation) the amount of training data required can easily exceed the maximum capacity of a single machine. [25] describes in detail how a new scientific field called *Distributed Learning* aims to distribute the training data and/or the neural network across many nodes, often implemented by combining an army of *slave* nodes that perform the calculations with a *master* node (often called the parameter server) that communicates with the slave nodes, iteratively combines their results, and updates the individual slave nodes with the combined result. This technique gained rapid popularity and is nowadays the backbone behind most industry-grade machine-learning implementations [26] (although there is also a multitude of other distributed learning architectures, each with their own advantages / disadvantages [25]).

Why federated learning?

Although distributed learning is highly effective in teaching neural networks to accomplish complex tasks, it still depends on as much data as possible to get the most accurate results. The data

to train popular neural networks often comes from smartphones, which on one hand produce enormous quantities of data which enables improved representation and generalization of machine-learning models, but on the other hand pose a significant problem because of three key reasons: (a) sending all kinds of data generated by smartphones over the internet consumes a lot of bandwidth, (b) training a neural network on data generated by billions of smartphones is computationally extremely intensive for a single *master* node, and (c) sending potentially sensitive information across the internet to the cloud raises privacy concerns, and is in certain cases not even allowed by several regulations such as the US HIPAA laws [27] and Europe’s GDPR law [28].”

These challenges motivated the development of a new type of distributed learning called *federated learning* where smartphones update the neural network with their data on-device and send back to the server the updated model rather than their data [29]. “Federated Learning brings the code to the data, instead of the data to the code” [30]. From this perspective, federated learning is closely related to Mobile Edge Computing (MEC) in the sense that computations are pushed to the edge of the network to reduce bandwidth consumption and improve privacy.

Thanks to these advantages, federated learning is used for a wide range of applications nowadays including next-word-prediction on keyboards such as Gboard [31-37], “wake word detection which enables voice assisting apps to detect wake word without risking exposure of sensitive user data” [38], speech recognition [39], wireless communications [40, 41], security applications (such as malware classification [42], human activity recognition [43], anomaly detection [44], and intrusion detection [45]), transport applications (such as data sharing between self-driving cars [46-48], preventing data leakage [49], traffic flow prediction [50], and the detection of attacks in aerial vehicles [51]), object detection [52], and health applications [53-57]

Federated learning environments have a number of notable characteristics [30, 58]:

- **Massively distributed:** the total number of nodes can easily be in the order of millions (**or even billions in the case of GBoard**).
- **Unbalanced:** the number of training samples per node can vary considerably.
- **Non-i.i.d. data:** different peers may possess different classes distributed in different ratios.
- **Unreliable:** since federated learning is often applied to smartphones, the nodes may go offline/online at any moment and the network connection may be slow.

However, we would like to build a system that not only seeks to properly handle the challenges imposed by the characteristics mentioned above but is also practical to be used in real-life. Therefore, we want to design the system around three additional principles:

Additional design principles

Decentralized

Practically all federated learning systems employ a centralized architecture characterized by a central trusted authority [30], often called a *parameter server*, that communicates with all nodes (generally smartphones). The machine learning process is as follows: (1) definition of the ML model (e.g. a CNN or RNN) by the developer in terms of hyperparameters, (2) distribution of the model by the master to the slaves, (3) local training of the model by each slave, (4) aggregation of all models by the master, (5) iteratively repeating steps 2, 3, and 4 [59]. The training process may stop when a sufficiently high accuracy is obtained or may be trained continuously as more data becomes available to the slaves.

Unfortunately, such a centralized architecture has several significant drawbacks [59-61]. The parameter server is not only a single point of failure susceptible to crashes or hacks, but it may also become a performance bottleneck when there are too many devices participating in the network. This problem accelerated research that aims to remove the parameter server entirely and train the network in a decentralized fashion [60, 62-64] where each node both trains and aggregates incoming parameters to learn the model [65]. Since these nodes are independent rather than being instructed by a master node, we won't refer to them as slaves but just as "nodes" or "peers" in the rest of the paper.

Byzantine-resilient

Even the most efficient and stable decentralized federated learning system is worthless for practical applications when the model can be ruined by Byzantine nodes, where *Byzantine* refers to the broadest class of faults in system components. The malicious model can be accidental (e.g. crashes, faulty sensors, computation errors, noisy transmission, nodes that are lagging behind, non-i.i.d. data, etc.; all of which the probability to occur at least once grows with the number of nodes [66]) or malicious on purpose (e.g. data poisoning or sophisticated model poisoning attacks; more on this in Section Data poisoning vs model poisoning attacks). This scenario, where the nodes don't know which of the other nodes are benign or corrupt, is the infamous "Byzantine Generals Problem" [67]. Without a Byzantine Fault Tolerance (BFT) mechanism, even a single malicious node that only takes moderate values to make its actions hard to detect can significantly degrade the performance of the federated model [68, 69].

Unfortunately, it's relatively easy to initiate a simple poisoning attack where a node just sends incorrectly updated parameters on purpose because of 3 major reasons [70]: (1) authentication mechanisms are often not feasible since federated systems often span across countries, (2) because the whole purpose of federated learning is to keep the training data private, it is impossible to audit the reliability of the training data, and (3) in real-world situations that dataset is often (very) non-i.i.d. which makes it challenging to distinguish between an attack and an unusual data class.

Byzantine resilience can be divided into weak and strong Byzantine resilience [71]. Weak f -Byzantine resilience implies

that despite the presence of f Byzantine nodes, the network will almost surely converge to a certain value. Strong f -Byzantine resilience implies that the network does not just converge in the presence of f Byzantine nodes, but also converges to approximately the right value. In this thesis we will focus on strong Byzantine resilience since this is the only provable solution against attackers with significant resources.

An interesting observation made by Haykin [72] is that a "mild" Byzantine worker can actually improve the performance of the system. This has to do with the fact that the optimization function of a neural network is practically never convex and has many local optima. By providing the "wrong" direction, a little bit of noise (or a "mild" Byzantine attack in that regard) can pull the optimization function out of a local minimum so that the network can converge to a better global minimum [73-75]. This is also the reason Stochastic Gradient Descent (SGD) works so well: a randomly drawn sample is inherently more noisy (higher variance) than the average of all samples [76] and may pull the network out of a local minimum. However, stronger Byzantine attacks can pull the network away from the global minimum in which case they ruin the network's performance.

Asynchronous

Distributed learning can happen either in synchronous rounds or in an asynchronous manner.

Although synchronous algorithms may seem to be the natural choice for federated learning (illustrated by the fact that the first federated learning protocol FedAvg [30] and influential subsequent research such as [77] were synchronous), there are a number of limitations as summed up by [78]:

- Some devices assumed to participate in the synchronous round may randomly drop-out, because of the volatile nature of the end-devices.
- Devices that just joined the network and are ready/willing to participate have to wait until the start of the next synchronization and are thus under-utilized.
- Some devices may be significantly slower than other devices due to more training data, less processing power, or less bandwidth.
- When a device is too slow to finish its iteration before the next synchronization, it has to overwrite its local model with the new global model and all of its progress is lost.

For this reason, numerous asynchronous approaches have been proposed [58, 59, 78-84]. These methods either overprovision clients and then accept the first x updates, dynamically update the synchronization time or amount of work per node, or use weighted averaging based on the staleness of incoming model updates. A common reoccurrence is that all of these approaches are dependent on a centralized parameter server, while in this thesis we aim to harness the advantages of a completely decentralized network. Truly decentralized protocols such as [85-87] are in practice always asynchronous because there is no server to synchronize training rounds, but these papers assume that the maximum staleness of peers is bounded.

[88] aims to achieve byzantine consensus in decentralized asynchronous networks, but they do not consider a situation where nodes can randomly join/exit the network.

[89] presents FedProx, a modification of the FedAvg algorithm that is supposed to tackle heterogeneity in FL by considering a variation of computational power and other differences between devices. However, its performance turns out to be sub-par in later research.

Another method based on FedAvg is SAFA [78], which aims to harness the potential efficiency gains of an asynchronous setting while using (a) a pace steering mechanism to reduce the impact of stale models and straggling clients, and (b) an aggregation algorithm that exploits a cache structure to reduce communication costs.

An approach that outperforms the former ones is presented in [90] that performs layer-wise matching and averaging of channels/neurons. It sends at the start of each training rounds global model matching results to the clients and adds additional neurons to the local models to achieve better performance.

To prevent the global model from drifting too much towards the fastest nodes, [91] proposes a mechanism to reduces this impact.

Key contributions

Overcoming all of the challenges described above is highly non-trivial, since distributing a computation over many peers induces a substantial risk of local crashes, computation errors, stalled processes, biased local datasets, but also possibly Byzantine workers trying to significantly degrade the performance of the system. Especially defending against Byzantine failures in a decentralized non-i.i.d. environment is challenging, while the literature on this topic is relatively sparse. Therefore, we will describe in this thesis Pro-Bristle, a novel algorithm that improves – to the best of the author’s knowledge – the current state-of-the-art by achieving (a) decentralized, (b) Byzantine-resilient, (c) asynchronous, and (d) non-i.i.d. federated learning. There are several types of machine learning, including supervised learning, unsupervised learning, and reinforcement learning. In this paper we will focus on arguably the most popular one, namely the supervised learning: a type of machine learning where the model should learn the correlation between the structure of the data and the corresponding label.

Overview of paper

First, we aim to give the reader an overview of the relatively unknown and emerging scientific field of Federated Learning. We will look at several types of Byzantine attacks, Byzantine-resilient defenses, and solutions to apply these techniques in a non-i.i.d. setting.

Then, we propose our solution called Pro-Bristle and explain how this gradient aggregation rule (GAR) improves the state-of-the-art in several domains.

Finally, the performance of this solution is compared with the performance of other GARs in a large variety of settings. We conclude by iterating gained and by giving directions for future research.

Related Research

Basic idea behind federated learning

In machine learning we aim to minimize the global cost function, risk function, loss function, or score $\ell(\theta)$ by finding the optimal model θ^* :

$$\theta^* = \underset{\theta}{\operatorname{arg\,min}} \mathbb{E}_{(x,y) \in D} \ell(f_{\theta}(x), y) \quad (1)$$

Where θ is the model, D is a distribution on $X \times Y$ and $\ell(\theta; i)$ is the loss of model θ on dataset instance i . This loss function is a proxy for the actual error to be minimized, generally the negative log likelihood of the ground truth class in the case of a classification problem.

This optimization problem is known as *risk minimization* but solving this problem is unfortunately for more complex models intractable. Therefore, a technique called Empirical Risk Minimization (ERM) is commonly used where take an empirically obtained dataset M i.i.d. sampled from D . Then we can obtain an estimate of the optimal model by calculating:

$$\theta_n = \underset{\theta}{\operatorname{arg\,min}} \frac{1}{|M|} \sum_{(x,y) \in M} \ell(f_{\theta}(x), y) \quad (2)$$

To minimize this function, a popular technique is called Gradient Descent (GD) that iteratively takes the derivative of the loss function with respect to the training samples and then moves the hyperparameters in the direction of the gradient:

$$\theta^{t+1} = \theta^t - \lambda \nabla_{\theta} \ell(\theta; i) \quad (3)$$

However, because the dataset may be large, gradient descent can take a long time to converge. A faster approach, used by almost all learning algorithms today, is to use Stochastic Gradient Descent (SGD) [92] where a subset (a *minibatch*) of the dataset is selected to update the parameters in a particular iteration [93, 94]. As a result, SGD produces faster but noisier updates than GD, but this noise is not necessarily a disadvantage because it also helps the algorithm to escape local minima. An important requirement for SGD to converge is that each minibatch is an unbiased sample of the actual distribution, which is typically achieved through uniform random sampling [95].

The most straight-forward way to apply stochastic gradient descent in a distributed or federated setting is to use a single master node (parameter server) that distributes and integrates the global model and a number of slave nodes that train the model that they obtained from the master node and send the result back to the master node [96, 97].

In distributed machine learning the most trivial implementation is called Bulk Synchronous Parallel [98] and in federated learning FedAvg [30]. FedAvg simply aggregates the models owned by the peers by coordinate-wise weighted averaging. It

was introduced by Google [43] and is still researched extensively from both an applied and theoretical perspective [59].

Byzantine attacks

FedAvg, as described in the previous section, takes the average of all models at every iteration. Obviously, when a single Byzantine node transmits a model with extremely low or high values, the average significantly changes, and the model become worthless. Such an attack is easy to detect, but there are many more sophisticated attacks that can, even with a single Byzantine attacker, considerably reduce the model’s performance and are much harder to detect [68, 99].

Byzantine attacks can be classified based on certain properties as either a data poisoning or a model poisoning attack, and as either an untargeted or a targeted attack.

Data poisoning vs model poisoning attacks

Data poisoning attacks such as [1-17, 100-103] try to degrade the performance of the learnt model to such a degree that the model becomes worthless by training the network with dirty samples. They were introduced by [14] to destroy support vector machines and later extended to many other ML algorithms including neural networks. Without appropriate Byzantine-resilient defense mechanisms, a malicious agent can relatively easily manipulate the global model. The best researched type of attack is a convergence prevention attack [18] where the attacker wants to prevent the network from converging and reduce its accuracy up to a point that the model might be utterly ineffective indiscriminately for testing examples. One might think that sending completely random numbers is an effective attack. However, because the mean of completely random numbers is 0, the network will still converge when the standard deviation isn’t too extreme {Muñoz-González, 2019 #310} (in fact, adding noise to the parameters is a popular method called *differential privacy* that is used to improve the user’s privacy {Dwork, 2008 #307} {Abadi, 2016 #308} {Wei, 2020 #309}). A scenario where a malicious agent injects malicious data into a benign client’s dataset (better known as a data injection attack) is also considered data poisoning. Another notable example of a data poisoning attack is a label-flip attack, where the labels of two or more classes are changed [2, 104] {Tolpegin, 2020 #311}.

While data poisoning attacks are based on the manipulation of training data, model poisoning attacks such as [1, 2, 4, 10, 18-23, 68, 69, 105, 106] manipulate the model’s parameters directly before sending it to other nodes. Consequently, every data poisoning attack can be imitated with a model poisoning attack [107], but model poisoning attacks give the attacker full control over every single parameter and can thus be much more effective as recent research has shown [13, 68, 99, 106]. They can even be used to replace the entire global model with a model of the attacker’s choice (model replacement attack), given a carefully chosen scaling factor [2, 68]. However, there are also simple model poisoning attack such as the Gaussian attack, where some

		Data poisoning	
[1-9]			[1, 2, 4, 9-17, 21-23]
	Untargeted		Targeted
	[1, 2, 4, 9-17]	Model poisoning	[1, 2, 4, 9-21]

of the gradient vectors are replaced by random vectors sampled from a Gaussian distribution with large variances.

Untargeted vs targeted attacks

Another way to classify Byzantine attacks as done by [107] is to group them into untargeted and targeted attacks. Whereas untargeted attacks such as [1-9, 68, 69, 100, 105, 106] aim to prevent convergence and reduce the global model’s accuracy [68, 69, 105, 106], targeted attacks such as [1, 2, 4, 9-23] aim to alter the model’s behavior in specific situations while keeping the total accuracy as high as possible to mislead Byzantine-defense mechanisms [2, 13]. Without proper defense mechanisms federated learning is susceptible to both untargeted and targeted attacks [104].

Targeted attacks are also sometimes called (semantic) backdoors, Trojan threats, or stealth attacks, where the backdoor can target either a single class (a label-flip attack) or a class of samples (e.g. an almost invisible attacker-chosen pattern of pixels, i.e. a trigger) causing an image to be classified incorrectly). A particularly effective attack is described by Bhagoji et al. [19] who use an alternating minimization strategy (alternating between training loss minimization and the boosting of updates for the malicious objective). A more sophisticated attack is proposed by Xie et al. in [21] who notes that all backdoor attacks until then used embeddings of the same global trigger pattern for all adversarial parties, called centralized backdoor attacks by the authors. They then propose distributed backdoor attacks (DBA) where the global trigger pattern is decomposed into local patterns and which is then embedded to different adversarial parties, thus making the attack harder to detect, easier to bypass robust aggregation rules, and more effective. In line with this contribution, [18] shows that targeted model poisoning attacks can become both significantly more effective and harder to detect when adversaries are able to collude.

As mentioned in Section Byzantine-resilient, a little random noise can actually improve the convergence of stochastic gradient descent. That might lead one to think that simply clearing away large deviations might be an effective defense mechanism. However, [18] the authors show that this assumption is incorrect and propose another powerful attack “capable of defeating all state-of-the-art defenses” based on injecting values

that are just within the perturbation range (the range of values that the Byzantine-defense mechanism allows).

Targeted attacks are hard to detect, because the accuracy of the model may not necessarily be impacted for any of the samples that any peer has available, but only for samples with, for example, a specific pattern that no-one except for the attacker knows about. More specifically, detecting backdoors in a model is an NP-hard problem, by a reduction from 3-SAT [108], and unlikely to be detected using gradient based techniques. [108] explains how it is relatively easy to develop a so-called *edge-case backdoor* which forces a model to consistently misclassify seemingly easy inputs that are unlikely to be part of the regular training data. Because these targeted model poisoning attacks only need to modify a small part of the model [107], they look quite similar as benign updates and require fewer adversaries than untargeted attacks, being already effective with even a single-shot attack under certain conditions [2].

Byzantine-resilient defenses

There are several types of Byzantine-resilient defense mechanisms that are in the literature often segmented into distance-based defenses (based on the calculation of some kind of distance between potential malicious attack vectors and some other vector(s), usually efficient but also vulnerable to elaborately designed Byzantine attacks [18, 106]) and performance-based defenses (based on testing the accuracy of a potentially malicious model on a small representative dataset, which is usually computationally quite intensive and clearly dependent on the availability of a test dataset but also usually quite effective) [109]. Another way of segmenting these algorithms is based on whether or not they are centralized (require a central parameter server) or decentralized, or by their degree of dimensional Byzantine resilience [109] (namely, the maximum number of tolerated Byzantine workers).

To structure the myriad of Byzantine defense mechanisms which we will call Gradient Aggregation Rules (GARs) from now on in accordance with the literature, we will use more fine-grained categories in this thesis, based on the fundamental principle that the algorithms employ.

Algorithm	Convergence Rate	Statistical Learning Rate	Condition on (M, b)
CM [15]	$\mathcal{O}(c')$	$\mathcal{O}\left(\frac{b}{M\sqrt{N}} + \frac{1}{\sqrt{MN}} + \frac{1}{N}\right)$	$M \geq 2b + 1$
CTM [15]	$\mathcal{O}(c')$	$\mathcal{O}\left(\frac{b}{M\sqrt{N}} + \frac{1}{\sqrt{MN}}\right)$	$M \geq 2b + 1$
GeoMed [16]	$\mathcal{O}(c')$	$\mathcal{O}\left(\frac{\sqrt{b}}{\sqrt{MN}}\right)$	$M \geq 2b + 1$
Krum [17]	N/A	N/A	$M \geq 2b + 3$
Multi-Krum [17]	N/A	N/A	$M \geq 2b + m + 2$
Bulyan [18]	N/A	N/A	$M \geq 4b + 3$
Zeno/Zeno++ [20], [21]	$\mathcal{O}(c') + \mathcal{O}(1)$	N/A	$M \geq b + 1$
RSA [23]	$\mathcal{O}\left(\frac{1}{T}\right) + \mathcal{O}(1)$	N/A	$M \geq b + 1$
signSGD [24]	-	N/A	$M \geq 2b + 1$

Distance-based screening

Screening potentially malicious incoming model updates for their distance with respect to the peer's own trusted model is by far the most popular to evade Byzantine attacks which should not come as a surprise. They are often highly efficient, do not depend on an extra dataset, special hardware features, or additional server, and they can do an excellent job to guard against

relatively simple attacks. However, although this class of algorithms is effective against simple attacks such as Gaussian and label-flip attacks, they are doing a poor with regard to more advanced attacks [110]. This is due to an implicit and somewhat incorrect assumption made by distance-based GARs, namely that close distances between model parameters implies similar performance. Additionally, gradient updates might differ significantly in a non-i.i.d. environment between nodes which result in large distances, leading distance-based GARs to reject these updates as outliers. Therefore, in the other sections other methods that are computationally much more expensive are evaluated that might be more effective than distance-based GARs, especially in non-i.i.d. settings.

Detecting outliers in non-distributed settings has been studied extensively for a long time [111], generally to be able to sanitize the data from poisoned or otherwise anomalous data [112]. In recent years, much progress has been made in terms of improved accuracy in high-dimensional settings [113-115]. For example, [17] uses a clustering technique to measure the difference between benign and malicious updates. However, these techniques are not suited for the distributed setting on which we are focusing.

A particularly influential algorithm is Krum[68] which selects the model that is most similar to all other models as the new global model. Even when the model being selected is malicious, the performance should not degrade too much in theory because it is close to all other models. Despite theoretical guarantees for the convergence for certain objective functions, Krum seems to perform quite bad in comparison to other algorithms [116] and often converges to an ineffective model [105]. The deficient performance stems from the ability of Byzantine workers to introduce a substantial change in a single parameter without significantly influencing the total distance due to the typically high dimensionality of the parameters [105]. [18] elaborates upon this insight and argues that, since only a single model is selected and even the best benign worker will have a few parameters that reside far from the mean, the GAR performs worse than other GARs that do integrate data from multiple models into the final model. The authors also briefly discuss Multi-Krum, which achieves similar accuracy at a faster rate by using an average of m local gradients obtained by Krum.

[99] presents two simple distance-based GARs, namely Coordinate-wise Trimmed Mean (CTM) (also evaluated by [117, 118]) which simply cuts off the smallest and largest b values in each dimension of the incoming vectors, and Coordinate-wise Median (CM) or Marginal Median which takes the median in every dimension. CM does not need at least $2b+1$ values like CTM, but does incur a performance hit because every dimension must be sorted to obtain the median.

[109] also evaluated CM and compared it under con-convex settings with two other approaches, namely the Geometric Median and Mean around the Median. The Geometric Median is defined as [68, 69, 109]:

$$\underset{v \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^n \|v - \tilde{v}_i\| \quad (4)$$

Which can be interpreted as the point for which the square distance to all other points in an n -dimensional space is minimized. The Mean around the Median is defined as the mean value of the $n - q$ indices closest to the median, where q is an arbitrary value. The authors find that Krum, Multi-Krum, and the Geometric Median perform worst, the Marginal Median has considerable variance, and the Mean around the Median performs best. The Geometric median not only performs poorly, but also dominates the training time in large-scale settings [119].

A variant on the *Geometric Median* is the (*Geometric*) *Median of Means* [69, 120-123], which first partitions all received vectors into k batches, then computes the mean for each batch, and finally takes the geometric median of the k batch means. [69] extends the techniques described in [124] with arbitrary/adversarial outliers. They only consider strongly convex losses which they try to remedy by using mini batches. However, their algorithm fails even when there is only a single Byzantine node in each mini-batch and is thus not very reliable.

Mhamdi et al. [105] try to combine the strengths of Krum and CM in Bulyan, a GAR that iteratively applies Krum to select a number of models and then a variant of CM to aggregate them into a single model. More specifically, Bulyan finds for every dimension the n parameters closest to the median and then takes their mean value. A notable disadvantage of Bulyan is its speed and the stringent condition that it imposes on the number of Byzantine nodes, namely $\#nodes \geq 4x \#nodes_{byzantine} + 3$. A year later, the authors extend Bulyan to Multi-Bulyan similarly to the extension of Krum to Multi-Krum, but unfortunately they don't report on the results [71].

In order to reduce the communication necessary for the aforementioned GARs, Bernstein et al. developed SignSGD [125] which only transmits the sign of every dimension of the gradient at every iteration. Since the global model is updated with an element-wise majority vote on the signs of the received gradients, the algorithm is in fact a median-based algorithm which makes it also robust against certain Byzantine failure and guarantees convergence given certain conditions imposed on the noise [126]. However, these conditions are typically not in order in a typical federated learning setting where the data is distributed non-i.i.d.[127]. Sohn et al. make SignSGD more robust against MITM-attacks, but do not address the case where nodes themselves are malicious[128].

In contrast to SignSGD, the work done by Li et al. [129], confusingly called RSA just like the cryptosystem, is able to handle heterogenous datasets in a Byzantine distributed setting and prevent incorrect gradient aggregation by letting every node store and update a local version of the global model which are then aggregated at the server by means of an ℓ_p -norm regularization term which regularizes the magnitudes of malicious messages.

An interesting distance-based method is described in [130] where the authors construct a graph where the nodes (representing models) are connected by a vertex only when their Euclidean distance is small enough and subsequently solve the maximum clique model to find the set of models that are similar to each other and therefore probably benign. Unfortunately, the authors only evaluate trivial label-flip attacks which make it hard to estimate the effectiveness of the algorithm in a more challenging environment.

All GARs described until now assume a federated setting where a single parameter server iteratively updates the global model. However, these algorithms do not translate well into a decentralized setting (which is the focus of this thesis) because decentralized GARs require consensus between all peers which is usually not required for distributed learning. To the best of our knowledge, there are only three papers that attempt to achieve Byzantine-resilient decentralized learning by adopting a truly distance-based strategy, namely ByRDIE [131], BRIDGE [132], and some extension of BRIDGE [133]. The CM algorithms described before (namely the ones presented in [99, 117, 118]) are suboptimal in vector-valued problems, because simply minimizing the objective function along each coordinate independent of all other coordinates yields the wrong solution (unless all dimensions are truly independent, which is generally not the case). This limitation is overcome in ByRdiE [131] by cyclically updating every coordinate one by one in a decentralized manner and subsequently applying trimmed-mean screening to obtain the final coordinate for each dimension. Given a strictly convex loss function, ByRDIE is proven to always converge (although not necessarily to an optimal solution). However, although ByRDIE might be efficient in terms of required training samples, ByRDIE is inefficient in terms of network communication because it only updates one coordinate at a time and the update step depends on the updates of other coordinates [116]. Therefore, [132] present BRIDGE which combines CTM with SGD (Stochastic Gradient Descent) to accomplish decentralized Byzantine-resilience with significantly less network communication for highly dimensional problems. The review paper [116] shows better performance for BRIDGE than for CTM, which is highly surprising because BRIDGE boils down to CTM in distributed environment. Upon closer examination this happens because the authors use a 0.7 connection ratio between the nodes to evaluate BRIDGE and just a $4x \#max_byzantine_nodes + 1 = 4x 2 + 1 = 9; 9 / 20 nodes = 0.45$ connection ratio between the nodes to evaluate CTM. [133] shows that the performance of BRIDGE

can be improved by adding a total variation (TV) norm penalty to allow some outliers to be able to handle non-i.i.d. data. This probably reduces the ability of the algorithm to defend against noise attacks, but unfortunately the authors have conveniently omitted these results from the Results section. [134] also builds upon BRIDGE and extends the solution to non-i.i.d. settings, but does this by re-introducing the central server that we want to omit in a decentralized setting (more information about non-i.i.d. approaches will be discussed in Section **Error! Reference source not found.**).

One of the most recent papers about the topic is [135] which select the models with the smallest Euclidean norm to be averaged for the updated model, but this paper evaluates its performance by measuring the loss function, which is extremely noisy and relatively unreliable compared to an evaluation using a validation dataset. This makes it hard to properly estimate its performance.

Performance screening

Although distance-based screening methods can be quick and effective to filter out “unusual” models, they will not include benign models when these models are quite different from the other models (e.g. in a non-i.i.d. environment) and an attacker can let the model drift towards a bad solution. Performance-based solution such as [12, 112, 136-139] detect malicious models based on a negative impact on the model’s accuracy given a test dataset. Of these papers, we want to seriously criticize the paper written by Zhao et al. [138] because, aside from the fact that it contains serious grammar errors and completely incorrect references, it also includes a major error about when label-flipping attacks are preferred above backdoor attacks. The authors say that label-flipping attacks are more effective in a scenario where data samples with the same label are quite similar while the latter is more suitable for scenarios where samples with the same label are quite diverse. This is incorrect: you want to use label-flipping attacks as an effective way to fool or prevent convergence of a model without any serious byzantine-resilient defense mechanisms while you want to use backdoor attacks to trick the model to misclassify certain input data without letting anyone notice that you are malicious (see Section Byzantine attacks). The authors also assume that agents share which labels they own, which is absurd: the whole purpose of a federated learning environment is that the user’s data (including the labels) stays private.

[136] presents RONI (Reject On Negative Influence, which removes training examples with a negative impact on the accuracy of the model) and [5] presents TRIM (which finds a subset of the training dataset given a pre-specified size and set of hyperparameters that maximizes the accuracy and is, according to the authors, more effective than RONI), both of which were originally intended to filter out bad training data on a single node. [106] converted and applied RONI and TRIM to a federated setting and found that in that setting RONI gives slightly better performance.

Xie et al. presented Zeno [140] (for synchronous environments) and Zeno++ [141] (for asynchronous environments), both of which use a centralized oracle that estimates based on a validation dataset the

true gradient and only keeps the k gradients most similar to this estimation. The performance of both GARs is quite good according to [116], but they depend on a centralized parameter server and need access to a sufficiently large unbiased validation dataset.

[142] takes a very different approach and proposes a GAR called PDGAN that uses a Generate Adversarial Network (GAN) to reconstruct the training data used by the peers to train the network. Based on this data, the accuracy of the received models can be estimated reliably. However, since the training data used by the peers is supposed to stay private, it is actually quite disturbing that GANs are able to reconstruct this data [19, 143], and are, in that regard, also a “highly impactful and prioritized” [144] attack in their own right.

Another recently presented GAR that is an important inspiration for this thesis is Mozi [110]. It first applies a distance-based strategy to quickly select a candidate pool of probably benign nodes, and then screens the resulting nodes based on their performance on a test dataset (performance screening).

Pruning

Since backdoor attacks (see Section Data poisoning vs model poisoning attacks) are extremely challenging to detect, let alone defend against, an entirely different class of defense mechanisms called “pruning” defenses has been proposed that specifically aims to prevent these backdoor attacks [145-147]. Pruning defenses use a representative subset of the global dataset (partially violating the FL assumption [107]) to evaluate which neurons in the neural network are inactive. These neurons are important to find and subsequently remove, because they enable attackers to create a backdoor in the model [14].

Unfortunately, even when these inactive neurons are removed from the model, more adaptive poisoning attacks are still possible [148]. After all, the boundary between a neuron being unused or being actively used is vague.

There are several other methods aimed at detecting backdoor [137, 145, 146, 149-154], but these methods either assume that there is a central server that can access the whole training dataset and scan the samples for malicious samples (which clearly is not possible in a federated learning setting) or access to a holdout set of similarly distributed data (which cannot help defend against more sophisticated model poisoning attacks as discussed in Section Byzantine attacks).

Behavioral-based

An original way to defend against targeted poisoning by sybil clones is presented in [104] and named FoolsGold. The authors first observe that, when sybils collude to poison a model, their “behavior is more similar to each other than the similarity observed amongst the honest clients”. Based on this insight, they present FoolsGold which detects and rejects poisoned contributions. However [138] shows that FoolsGold is unable to recognize against a powerful single-client attack and can be evaded by decomposing a distributed attack into several orthogonal vectors.

Whereas all GARs discussed until now make it as difficult as possible for an attacker to manipulate the system, there is also a wide variety of GARs that take a different approach and aim to eliminate any incentive for a node to attack the system. A trivial approach where a parameter server simply assigns a reputation based on a performance-based screening procedure per node (such as [155]) doesn't work well, because a malicious attacker can first build up an excellent reputation, and then suddenly completely ruin the model, empowered to do so thanks to its good reputation. A better approach turns out to be rewarding and punishing of participants based on their contributions, something that can be facilitated in decentralized environments by means of a distributed ledger [46, 156-170], usually a blockchain. This ledger can also be used to save global model parameters to enhance the security of the system [161, 166].

Kang et al. introduced in [171] reputation as a means to determine the reliability of every node and subsequently proposed a GAR based on these reliability scores [172], using RONI to calculate reputations in i.i.d. environments and FoolsGold to calculate reputations in non-i.i.d. environments. For this to work, the authors (implicitly) assume an environment where nodes have a strong identity and where there are many different parameter servers learning different tasks that share reputation opinions of nodes over a public blockchain.

[173] also assign a reputation to nodes that contribute well, but their algorithm is seriously flawed: the authors use KRUM to determine if an update is benign (which is highly unreliable [116]) and then increase/decrease a node's reputation when the update is accepted/rejected respectively, implying that you can make theoretically for every good contribution also a bad contribution (in practice a single bad contribution can damage the model significantly while the impact of a single good contribution is generally very limited).

Whereas all former approaches assume that the individual workers might be Byzantine, [174] assumes a centralized setting where the parameter server might be Byzantine. They use a blockchain to audit all model updates from all peers so that everyone can verify if the parameter server aggregates the model updates correctly. The authors also train an autoencoder to recognize outliers (i.e. Byzantine attack). This seems to work quite well, but the autoencoder is only effective after it has been trained properly which may take many iterations.

Another blockchain-based approach is described by [156] called HoldOut SGD which first splits the nodes into a set of workers that use their data to train the model for a single iteration, and a voting committee that votes for the best proposals and stores this information on the blockchain (similarly to [159, 175]). The voting committee is usually selected based on Proof-of-Stake and Verifiable Random Functions (VRFs)). The method is fully decentralized but defends only against up to a factor of $1/3^{\text{rd}}$ Byzantine workers. The technique is hardly scalable to a high number of nodes, because every node in the voting committee has to evaluate every single update, and because either all voting committee nodes are waiting until the workers are finished or vice versa a significant amount of time is spent idling for every node.

Where the blockchain papers mentioned above are of good quality, one has to be very careful when search for literature about this subject. There are many papers where a blockchain is used without properly understanding its (dis)advantages. For example, in [46] the authors say that they want to address privacy issues by using a blockchain, but simply using a blockchain doesn't magically improve the user's privacy. The authors also states that Directed-Acyclic-Graphs (DAGs) are a certain kind of blockchain (which is incorrect, it are different technologies. Stating that DAGs and blockchains are both examples of Distributed Ledger Technology (DLT) would have been correct) and that DAGs use cumulative PoW, which is also incorrect: newly added data DAGs usually just reference and validate previous transactions without any PoW involved.

A particularly good paper where the authors really make use of the strengths of DLT is [157] where they use a Tangle to represent the approved transactions as nodes in a Directed Acyclic Graph (DAG). Every new transaction first verifies two previous transactions by using one of the defense mechanisms described in this section and includes the updated model parameters.

There is also a considerable body of literature that uses behavioral techniques to incentive nodes with high quality training data to participate in the training process such as [163, 165, 168-172, 176-188] and Stackelberg game methods [176, 177, 189, 190], but since these methods are not supposed to defend against Byzantine attacks, we leave them out of the scope of this thesis. Additionally, the underlying assumption, namely that agents should get some kind of incentive to participate in a federated training process, does not seem to apply in many popular applications, such as Gboard, Captcha, or Google Fit.

Other

There are several papers discussing innovative defense mechanisms that are not easy to classify into a particular category. There are a few papers that use Trusted Execution Environments (TEEs) to achieve some form of security such as [163, 191-194]. [191] uses secure aggregation based on the Secure Multiparty Computation (SMC) algorithm to aggregate the values of untrusted nodes without revealing these values, enabling a parameter server that every party can trust. [192] also discusses how TEEs can be used as a defense technique, later used by [193, 194] to create a generic framework that can be used to integrate TEEs in a federated learning environment. Weng et al. [163] developed DeepChain that on one hand uses a blockchain to incentivize parties to participate in the training process, and on the other hand uses a combination of Intel Software Guard Extensions (SGX) enclaves and homomorphic cryptographic functions to provide a safe and privacy-preserving environment. Their solution work well, but is computationally also very expensive, limiting its use cases.

There are also a few solutions that model defending against Byzantine attacks as a learning problem. [195] uses a Recurrent Neural Network (RNN) and an auxiliary dataset to aggregate gradients in a Byzantine-resilient manner. The idea is that a

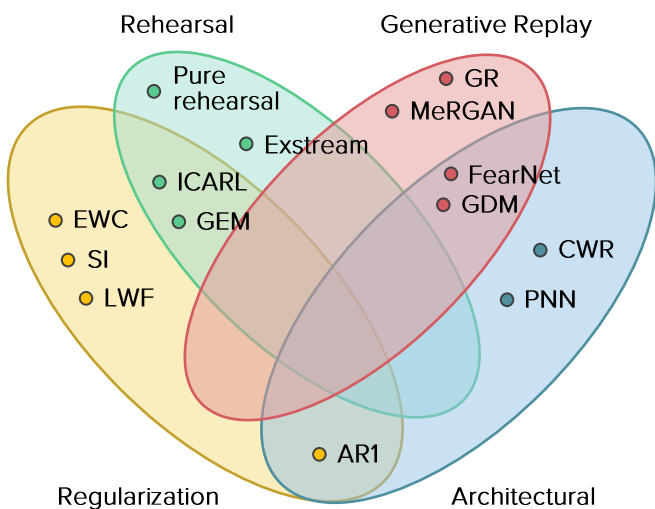
machine learning approach can detect attacks that is hard for other algorithms to accomplish. Unfortunately, since their RNN is a “black box”, the authors are unable to give any theoretical guarantees. [196] uses variational autoencoders with spectral anomaly detection to detect malicious updates based on their low-dimensional embeddings. By removing the noisy and irrelevant features, the anomalous (malicious) model updates can be distinguished from the benign updates in a low-dimensional latent feature space.

A final type of defense mechanisms we would like to highlight are defenses based on replicating the same training over multiple nodes [119, 128, 197, 198]. When all nodes are benign, they will obviously report the exact same results. While the accuracy of these mechanisms is often illustrated with rigorous theoretical guarantees, they generally assume a centralized server with either a copy of the data or the ability to globally shuffle the data, which makes the algorithm inappropriate for a decentralized federated environment. For example, DRACO [119] lets the parameter server send to multiple workers the same chunk of data and uses majority voting to find the correct evaluation. Against a small number of Byzantine attackers DRACO is very robust, but the algorithm scales poorly to a greater number of attackers. For example, when there are just 5 attackers, each chunk already needs to be calculate $5 \times 2 + 1 = 11$ times.

Non-i.i.d.

Whereas in regular distributed learning environments, a characteristic of a typical federated learning environment is that the shards of the slaves are non-i.i.d. (not independent and identically distributed) [68, 69, 99]. For example, it is possible that device A has class X and Y, and device B has class Y and Z. As a result, the model of device B will be quite different from the model of device A, making it hard for device A to determine if device B’s model is malicious or not. To make matters worse, a trivial average of the parameter updates yields a considerably worse model than a model that would have been trained by a single node on class X, Y, and Z [30, 199, 200].

The challenge of building a single global model by combining



multiple local models without reducing their accuracy is closely

related to *multi-task learning*[201]. Multi-task Learning or Continual Learning (CL) is concerned with preventing *Catastrophic Forgetting* or *Catastrophic Inference*, a phenomenon where the neural network completely forgets what it has learnt before when it is taught a new task. Instead, the network should be capable of *Lifelong learning*: continuously acquire new knowledge, refine existing knowledge, and prevent new tasks from interfering with existing knowledge.

Figure 1 is a Venn diagram created by Lesort and Lomonaco [24] categorizing the existing CL methods into 4 partially overlapping categories:

Architectural approaches seek to allocate additional neural nodes whenever they are required or freeze specific weights[202-205], but this requires the developer to know the number of tasks / samples per task a-priori and leads to scalability issues for large neural networks.

Regularization techniques minimize the extent to which the most important weights are overwritten by the training on a new model. Elastic Weight Consolidation (EWC) [205], which was based on Learning without Forgetting (LwF) [206] is an influential regularization technique that extends the loss function with a quadratic penalty on the change in parameters that are important for the formerly learned tasks. The authors set the importance of the parameters to the diagonal of the Fisher information matrix, which works well for learning permutations of the same task, but not for learning entirely new categories

FIGURE 1. Venn diagram of existing CL methods[24]

incrementally [207]. Several improvements have been made since such as [208-211].

Rehearsing old samples interleaved with new samples is also an effective way to prevent catastrophic forgetting. [200] concluded that globally sharing just 5% of the training data can result in a 30% increase in accuracy. These training samples can be selected randomly or carefully to be as representative of the coresets as possible. However, this approach increases the amount of memory needed to store all samples[212-215].

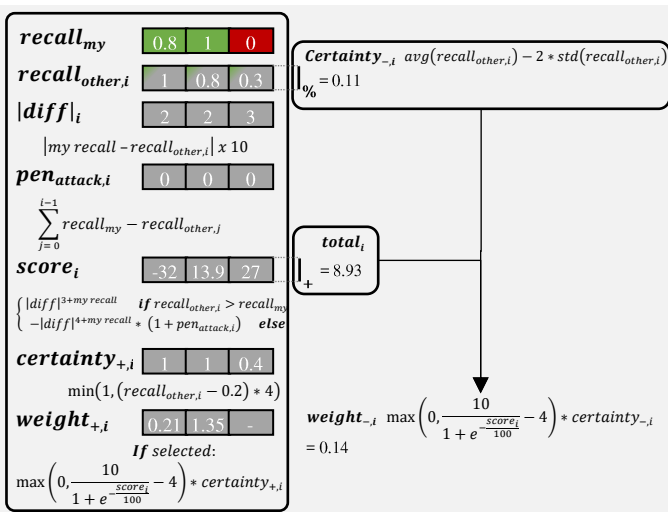
Generative replay is a variant on rehearsing old samples where a *Generative Adversarial Network* (GAN) is used to artificially generate samples that have a similar distribution as the past experiences. These samples are then intertwined with the new empirical training samples just like in rehearsal-based strategies.

The approaches discussed until now are generic multi-task learning techniques, but there has also been research into similar techniques specifically for federated learning environments.

[29] is an example of a non-i.i.d. approach for federated learning, but the authors use clusters which do not work well on high-dimensional data (such as neural networks): the authors simply throw away all parameters of the neural network except for the first 288 parameters in the first layer. The technique presented by [200] is more effective and uses rehearsal: they assume that a small amount of i.i.d. data is available that can be shared across all peer nodes (which is generally a realistic assumption).

In specific situations where the loss function is convex and its conjugate dual is expressible, research has shown that dual coordinate ascent approaches such as Mocha en Cocoa can yield superior results [200, 216-218]. Mocha [218] for instance handles non-i.i.d. datasets while also tackling the challenge of fault tolerance, stragglers, and communication efficiency. The algorithm models the relation between the tasks by adding a loss term and subsequently uses a primal-dual formulation to solve the optimization problem. However, like many other multi-task learning algorithms, it assumes that all peers participate in every training round which makes these algorithms harder to apply in a truly federated setting.

A particularly popular approach seems to be to use Elastic Weight Consolidation ([219-223] which, as explained in this section before, penalizes large changes of parameters important for previously learned tasks. It is somewhat surprising that more recent methods such as CWR(+), LWF, or AR1 have not been investigated yet because these methods perform significantly better in non-federated environments than EWC[224].



Proposed solution: Pro-Bristle

The main contribution of this work is the development of a new GAR named Pro-Bristle (Practical yet RObust Byzantine-Resilient decentralIzed StochasTic federated LEarning). To explain how Pro-Bristle works we will first recapitulate the four characteristics and three additional design principles that we listed in Section Introduction).

A federated learning environment is characterized several characteristics:

- Massively distributed network
- Unbalanced data
- Unreliable nodes
- Non-i.i.d. data

Additional design principles our GAR will obey are:

- Decentralized
- Byzantine-resilient
- Asynchronous
- Communication-efficient

In contrast to practically all existing works on federated learning until now we will use *gossiping* to make a **massively distributed** number of nodes learn together in a **decentralized** and scalable way. Gossiping also makes the fact that the nodes are **unreliable** irrelevant, since gossiping happens with “a random node”; if some node happens to be offline, the other nodes will just choose other nodes to gossip with. To provide **Byzantine-resilience** we first observe that in a perfect world (non-i.i.d., all peers working synchronously on the same iteration, and with a balanced dataset) for every node A , all benign models that node A receives will be reasonably close to node A 's own model. However, Byzantine models that the node receives can be either within or outside this distance. This inspires us to first apply a distance-based filter to get rid of some Byzantine attacks, and then apply a performance-based filter to filter out more sophisticated Byzantine attacks.

However, in real-life, the world is not perfect (as illustrated in Section Why federated learning?). For example, in a gossiping decentralized environment, it is natural for the nodes to operate in an **asynchronous** manner which makes things more complicated. In an asynchronous environment we will differentiate between three problems:

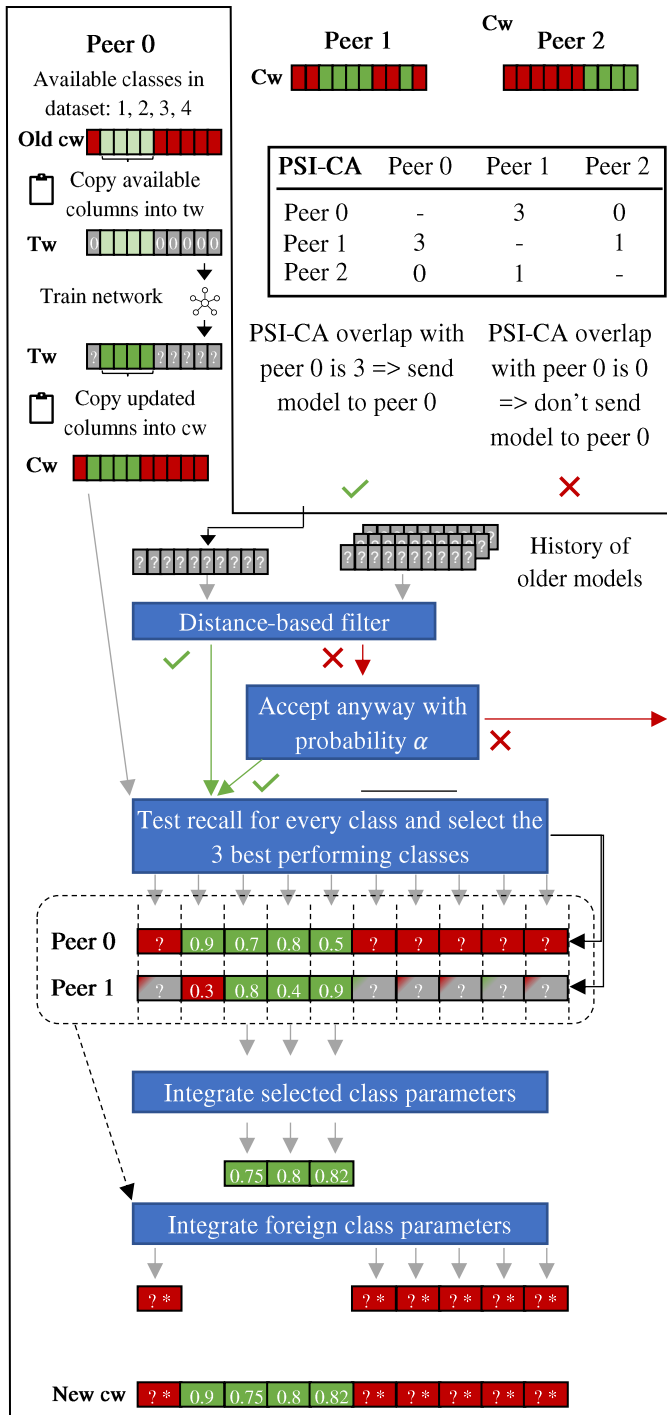
- **Too few peer models received to reliably perform distance-based screening.** After an iteration (which is typically significantly faster on an asynchronously operating node since the node does not have to synchronize with (=wait for) all other nodes) the number of models received from other nodes is arbitrary. When iterations are computed fast (e.g. because the neural network is small) and models are received only slowly (e.g. because the bandwidth is limited), the number of received models might be small. When, for example, only two models are received, trivial distance-based screening procedures obviously do not work because there are insufficient nodes to compare the models received with. We propose to add a buffer to keep track of a number of recently received models

to make the distance-based screening procedure more robust. Applying this idea in a federated setting is not entirely new because it was also explored by Yang et al. [83]. However, the paper of Yang et al. (a) assumes a centralized instead of a decentralized setting, (b) does not explicate what advantage using multiple buffers exactly gives in their solution (which is entirely unclear to us), and (c) is unable to update its model directly after receiving an update, resulting in subsequent local training on a (slightly) outdated model.

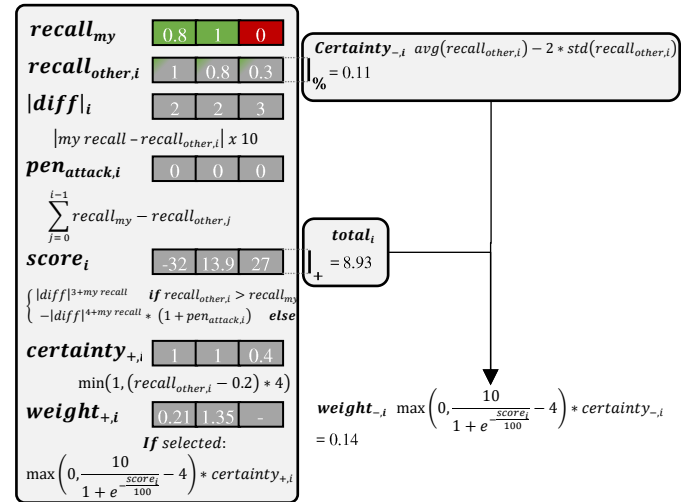
- **Stale model received.** A model that is received might be outdated and stale. In this case, its performance will be subpar and therefore be filtered out either in the distance-based screening or in the performance-screening phase.
- **Extremely accurate model received that is so different from the peer's own model that it is filtered out by the distance-based screening method.** A model that is received might be way better than the current model, causing it to be filtered out by the distance-based screening method. To solve this, we propose to use a popular strategy that has never been applied in a federated learning setting before: exploration vs exploitation. Based on an exploration ratio α , the distance-based screening filter should randomly accept models that are “not close enough” to be considered otherwise. The performance-screening filter will then notice the supposedly superior performance of this model. We also propose that to use weighted averaging to aggregate the models based on their performance: when a model is received that performs extremely well, the node should shift its current model very significantly towards this model.

Finally, to properly function in a **non-i.i.d.** environment, we need to address two challenges:

- We need to be able to prevent Byzantine attacks, which is non-trivial because the model of a peer with completely different data will likely be recognized as malicious by both the distance-based and performance-based screening procedure. To solve this, we first take a step back and determine if a node has enough overlap in its data distribution with another node to properly check the accuracy of the received model. Unfortunately, in a federated setting it is not possible to simply compare the data of a pair of nodes, because this data is private. Therefore, we use Private Set Intersection Cardinality (PSI-CA) to check the overlap between the datasets. Only when nodes have sufficient overlap, they will gossip with each other. We assume that there is sufficient overlap for the whole network to be well-connected.
- We need to be able to prevent catastrophic forgetting when combining models trained on different classes. Catastrophic forgetting means that when a neural network is trained for a certain set of classes and



thereafter is trained for another set of classes, it completely overwrites (forgets) the first set of classes. In a federated environment this results in nodes constantly overwriting each other, resulting in mediocre performance. Pro-Bristle is hard to compare to existing methods. For instance, it makes the (realistic) assumption that a large public dataset with approximately the same low-level features is available, which is used to pre-initialize the neural network.



Pseudocode

Input: initial estimate x_0 , dataset D^{train} containing an arbitrary non-validation subset of the node's collected data, dataset D^{test} containing a small trusted set of samples for each label, history buffer size γ , exploration ratio $\frac{\beta}{\alpha}$, max weight Ω , weight decay η , transfer network Ψ

$\mathcal{N}^s \leftarrow \text{getSimilarPeers}()$

$l \leftarrow$ initialize loss function by deep transfer from Ψ

For $t = 0, 1, 2, \dots$ **do**

Stochastically sample $\xi_i(t)$ from D_i^{train}

$\nabla l(\mathbf{x}_i(t), \xi_i(t)) \leftarrow$ Compute the local gradient

$\mathcal{N}_i^n(t) \leftarrow$ All models received from peers $j \in \mathcal{N}^s$

$\mathcal{N}_i^r(t) \leftarrow$ The γ most recent models received in previous iterations

if $|\mathcal{N}_i^n(t)| > 0$ **then**

For j **in** $(\mathcal{N}_i^n(t) \cup \mathcal{N}_i^r(t))$ **do**

$d_{i,j} \leftarrow \|\mathbf{x}_i(t) - \mathbf{x}_j(t)\|$

End for

$\mathcal{N}_i^d(t) \leftarrow \left(\underset{|\mathcal{N}^*|=\alpha}{\text{argmin}} \sum_{j \in \mathcal{N}^*} d_{i,j} \right) \setminus \mathcal{N}_i^r(t)$

$\mathcal{N}_i^e(t) \leftarrow \mathcal{N}^* \underset{\subseteq}{\text{N}^* \text{ is a random subset where } |\mathcal{N}^*| = \beta} (\mathcal{N}_i^n(t) \setminus \mathcal{N}_i^d(t))$

$\mathcal{N}_i^c(t) \leftarrow \mathcal{N}_i^d(t) \cup \mathcal{N}_i^e(t)$

For $j \in i \cup \mathcal{N}_i^c(t)$ **do**

For c **in** $\text{classes}(D_i)$

$\text{recall}_{j,c}(t) \leftarrow \text{recall}(l, \mathbf{x}_j(t), D_i^{test})$

End for

if $j \neq i$ **then**

$C_j(t) \leftarrow \underset{|\mathcal{C}|=\text{cardinality}(\mathcal{N}_j^s(t))}{\text{argmin}} \sum_{c \in \text{classes}(D_i)} \text{recall}_{j,c}(t)$

For c **in** $C_j(t)$

$\text{weightdiff}_{j,c}(t) \leftarrow |\text{recall}_{j,c}(t) - \text{recall}_{i,c}(t)|$

$\text{seqAttackPenalty}_{j,c}(t) \leftarrow \text{getSeqAttackPenalty}(C_j(t), c, \text{recall})$

if $\text{recall}_{j,c}(t) > \text{recall}_{i,c}(t)$

$\text{score}_{j,c}(t) \leftarrow \text{weightdiff}_{j,c}(t)^{3+\text{recall}_{i,c}}$

Else

$\text{score}_{j,c}(t) \leftarrow -\text{weightdiff}_{j,c}(t)^{4+\text{recall}_{i,c}} * (1 + \text{seqAttackPenalty}_{j,c}(t))$

End if

$\text{certainty}_{j,c}(t) \leftarrow \text{clamp}(0, (\text{recall}_{j,c} - 0.2) * 4, 1)$

Input: initial estimate \mathbf{x}_0 , dataset D^{train} containing an arbitrary non-validation subset of the node's collected data, dataset D^{test} containing a small trusted set of samples for each label, history buffer size γ , exploration ratio $\frac{\beta}{\alpha}$, max weight Ω , weight decay η , transfer network Ψ

$\mathcal{N}^s \leftarrow \text{SimilarPeers}()$

$l \leftarrow$ initialize loss function by deep transfer from Ψ

For $t = 0, 1, 2, \dots$ **do**

Stochastically sample $\xi_i(t)$ from D_i^{train}

$\nabla l(\mathbf{x}_i(t), \xi_i(t)) \leftarrow$ Compute the local gradient

$\mathcal{N}_i^n(t) \leftarrow$ All models received from peers $j \in \mathcal{N}^s$

$\mathcal{N}_i^r(t) \leftarrow$ The γ most recent models received in previous iterations

If $|\mathcal{N}_i^n(t)| > 0$ **then**

For j in $(\mathcal{N}_i^n(t) \cup \mathcal{N}_i^r(t))$ **do**

$d_{i,j} \leftarrow \|\mathbf{x}_i(t) - \mathbf{x}_j(t)\|$

End for

$\mathcal{N}_i^d(t) \leftarrow \left(\underset{|\mathcal{N}^*|=\alpha}{\operatorname{argmin}}_{\mathcal{N}^* \subseteq \mathcal{N}_i^n(t)} \sum_{j \in \mathcal{N}^*} d_{i,j} \right) \setminus \mathcal{N}_i^r(t)$

$\mathcal{N}_i^e(t) \leftarrow \mathcal{N}^* \underset{\subseteq}{\operatorname{argmin}} \mathcal{N}^*$ is a random subset where $|\mathcal{N}^*| = \beta \left(|\mathcal{N}_i^n(t) \setminus \mathcal{N}_i^d(t)| \right)$

$\mathcal{N}_i^c(t) \leftarrow \mathcal{N}_i^d(t) \cup \mathcal{N}_i^e(t)$

For $j \in i \cup \mathcal{N}_i^c(t)$ **do**

For c in $\text{classes}(D_i)$

$\text{recall}_{j,c}(t) \leftarrow \text{recall}(l, \mathbf{x}_j(t), D_i^{test})$

End for

If $j \neq i$ **then**

$C_j(t) \leftarrow \underset{|\mathcal{C}|=\text{cardinality}(\mathcal{N}_j^s(t))}{\operatorname{argmin}}_{c \text{ in classes}(D_i)} \text{recall}_{j,c}(t)$

For c in $C_j(t)$

$\text{weightdiff}_{j,c}(t) \leftarrow |\text{recall}_{j,c}(t) - \text{recall}_{i,c}(t)|$

$\text{seqAttackPenalty}_{j,c}(t) \leftarrow \text{getSeqAttackPenalty}(C_j(t), c, \text{recall})$

If $\text{recall}_{j,c}(t) > \text{recall}_{i,c}(t)$

$\text{score}_{j,c}(t) \leftarrow \text{weightdiff}_{j,c}(t)^{3+\text{recall}_{i,c}}$

Else

$\text{score}_{j,c}(t) \leftarrow -\text{weightdiff}_{j,c}(t)^{4+\text{recall}_{i,c}} * (1 + \text{seqAttackPenalty}_{j,c}(t))$

End if

$\text{certainty}_{j,c}(t) \leftarrow \text{clamp}(0, (\text{recall}_{j,c} - 0.2) * 4, 1)$

Implementation

Datasets

From all papers that we read as part of the literature review for this thesis, two datasets turned out to be extremely popular in the literature, namely the **MNIST** dataset [4, 11, 14, 17, 18, 20, 21, 29, 81, 89-91, 99, 100, 104-106, 108, 109, 116, 119, 125, 127, 130-135, 138, 139, 142, 143, 146, 148, 149, 151, 153, 155, 156, 159, 160, 162, 163, 173, 193-196, 199, 200, 225] and the **CIFAR-10** dataset [2, 16, 18, 21, 79, 81-84, 86, 105, 108, 109, 119, 125, 126, 128, 135, 137, 138, 140, 141, 147, 153, 156, 193, 195, 197, 199, 200, 220].

The MNIST dataset consists of 60,000 gray-scale training images and 10,000 test images of 28x28 px that represent handwritten digits. Even though MNIST does not represent a typical federated learning dataset, it is very popular since it is easy and fast for neural networks to learn, and, thanks to its status as one of the most popular machine learning datasets, makes it possible to compare our results with a large body of established literature.

The CIFAR-10 dataset also consists of 60,000 training images and 10,000 test images. These images are 32x32 px and RGB-colored, showing pictures of ten distinct types of objects such as cars, airplanes, and dogs. CIFAR-10 turns out to be significantly more challenging to learn than MNIST, which might be useful to properly investigate the power of new algorithms.

We also include a realistic federated learning dataset, namely the **UCI-HAR** dataset, one of the most popular smartphone datasets [226]. This dataset consists of 10299 recordings of people performing one of the six included activities. Every recording consists of 561 measurements of six coordinates measured by the gyroscope and acceleration sensors.

Machine Learning part

For MNIST and CIFAR-10, we use the same CNN architectures as used by [30] with the only difference that we use Leaky ReLU instead of the regular ReLU as activation function for the hidden layers, since the former one suffers less from the vanishing gradients problem. For the output function we use the softmax function and as loss function, we use negative loglikelihood.

MNIST

Layer	Details
Convolution	Kernel: <5, 5>, stride: <1, 1>
Max pooling	Kernel: <2, 2>, stride: <2, 2>
Convolution	Kernel: <5, 5>, stride: <1, 1>
Max pooling	Kernel: <2, 2>, stride: <2, 2>
Dense	#nodes: 500
Output	#nodes: 10

CIFAR

Layer	Details
Convolution	Kernel: <3, 3>, stride: <1, 1>
Batch normalization	
Max pooling	Kernel: <2, 2>, stride: <2, 2>
Convolution	Kernel: <2, 2>, stride: <1, 1>
Batch normalization	
Max pooling	Kernel: <3, 3>, stride: <1, 1>
Convolution	Kernel: <2, 2>, stride: <1, 1>
Batch normalization	
Max pooling	Kernel: <2, 2>, stride: <2, 2>
Convolution	Kernel: <2, 2>, stride: <1, 1>
Batch normalization	
Max pooling	Kernel: <2, 2>, stride: <2, 2>
Output	#nodes: 10, dropout: 0.8

HAR

Layer	Details
1D convolution	Kernel: <3>, #nodes: 64
1D max pooling	Kernel: <2>, stride: <2>
1D convolution	Kernel: <3>, #nodes: 64
Global max pooling	
Output	#nodes: 10

Gradient Aggregation Rules

To properly compare our proposed solution with existing methods, we implemented five other gradient aggregation rules (described in detail in Section Byzantine-resilient defenses):

- **FedAvg** [30]. FedAvg is equivalent to simple averaging, is researched extensively, and very often used as baseline to compare other GARs against.
- **CM (Coordinate-wise Median)** [99]. CM is perhaps the simplest, but also a very effective Byzantine-resilient defense mechanism, as shown by [116].
- **Krum** [68]. Krum is an extremely popular GAR that selects the model with the minimal local sum of Euclidean distances.
- **Bridge** [132]. A very recent survey paper [116], published in May 2020, concluded that Bridge was the best performing GAR in decentralized settings.
- **MOZI** [110]. MOZI was published shortly after [116]'s survey and uses a hybrid between distance-based and performance-screening to achieve superior results.

Environment

We use two separate ways to test the performance of Pro-Bristle, namely in a local simulation and in a truly decentralized environment. In the former case, we run a single program that iteratively trains and combines up to 250 models to simulate a small-scale federated setting. This approach is not only very fast, but also makes it easy to accurately control a variety of settings, such as little computation power, low bandwidth, nodes that randomly join/exit, etc. We also emulate 16 completely independent smartphones to test if the results are comparable in a "real" setting. Unfortunately, this limit of 16 emulators is hardcoded in the Android emulator executable which makes it

hard to run more emulators simultaneously. However, 16 emulators is enough to accurately measure the performance of different GARs and, if the programs works well on 16 emulators, gives us confidence that the code works as intended and will also scale to a higher number of nodes.

Network protocol

The nodes that use federated learning to collaboratively learn a model communicate with each other over a network. Since we aim to run everything completely decentralized, it is non-trivial for nodes to find and communicate with each other in a fault-tolerant and effective way. Therefore, we use IPv8[227, 228], a well-established decentralized peer-to-peer (P2P) middleware stack that is used by i.a. the popular Tribler media sharing system[229, 230]. Furthermore, we extended IPv8 with two significant performance improvements to make the system more effective.

The first improvement is an extension to the Trivial File Transfer Protocol (TFTP) that enables parallel transmission of multiple files between the same 2 nodes. This was implemented by assigning a unique file identifier to each file and prepending every data packet with this identifier to keep track of all packets. The second improvement is, to speed up the slow transmission times of TFTP, the first Kotlin implementation of the micro Transport Protocol (μ TP). This protocol aims to mitigate the poor latency and congestion control problems found in regular TCP implementations, while providing reliable and ordered packet delivery. It sends multiple packets simultaneously and automatically slows down the transmission when the network seems to get congested.

Task automation

Creating and starting all emulators, and installing, starting, initializing, running, and evaluating the federated learning program on every emulator, is infeasible to do by hand for a large number of emulators. Therefore, we created a separate coordinator that automates these tasks. Based on the current operating system it executes several scripts (for example to create new emulators that are reset to factory settings, or to increase the local network buffers to decrease the uncontrolled/unintended packet loss to speed up the network communication) to create and run all tests consecutively. The nodes are instructed to communicate their evaluations to the coordinator, who writes the evaluations to a CSV file.

A dedicated Python script was used to process the evaluations and generate the figures as shown in this paper.

Biggest issues encountered

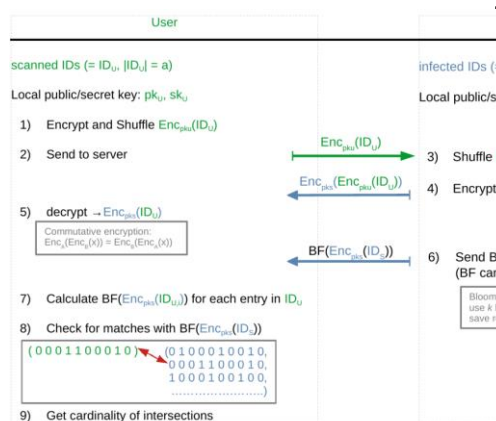
- Weighted averaging issues
 - o Problem with communicating new parameters
 - o Problem with communicating gradients
- Local network buffers
- Elastic weight consolidation
- Limit of 16 emulators
- TFTP insufficient
- Debugging UTP
- DL4J dependency hell
 - o NaNs door verkeerde versie
 - o Verouderde URL gehardcoded
 - o Not maintained anymore
- Multithreading issues
 - o Had to use ConcurrentHashMap for connectionIds
- Sometimes sending a message to the other peer and receiving the subsequent response happened faster than executing the next line of code
- DL4J gives slightly different values when adding and then subtracting instead of subtracting and then adding
- DL4J has a bug in its memory management when it runs multiple threads simultaneously => solved by running them sequentially

Threat Model

- Specify threat model (e.g. like in [106])
- Important insight: all benign models are within a reasonable distance of the node's own values; byzantine models can be within or can be outside this distance => distance-based filter as first rule, and only then a performance-based filter
- Specify problem formulation like in [139]
- Often modeled as Poisson process [139, 231, 232]
- [233, 234]
 - o Attacks that break existing defenses against Byzantine adversaries
- Assume that each node is connected to at least 1 benign node
- [139]
- Checks if incoming gradients are similar to own gradients => no updates possible anymore when model becomes stale
- Many approaches assume that the number of adversarial workers is less than half of the total number of workers, are some exceptions that ensure convergence even in the presence of a large number of adversaries, namely [139, 141, 155, 195, 235]
- "we assume a network population with hundreds or thousands of devices that are not typically available at the same time to perform training; furthermore, limitations in compute and storage resources, as well as network bandwidth, are to be expected [21]. The training dataset is horizontally partitioned, i.e. devices have different sets

of non-IID training and validation examples that include a common set of features."

- "We are focused on FL and therefore assume that data is distributed across clients and hidden, such as in an IoT deployment with multiple devices distributed in people's homes. The adversary can only access and influence the model state through the FL API. They cannot observe the training data of other honest clients. The adversary can observe the global change in model state to learn the total averaged update across all clients, but they cannot view individual honest client updates."
- [236]
- Use averaging across iterations, but critically depends on a parameter server that keeps track of all gradient updates of all nodes; also "proposes a fault-tolerant SGD variant different from the robust aggregation rules. The algorithm utilizes historical information, and achieves the optimal sample complexity."
- "Alistarh et al. propose a Byzantine-resilient SGD algorithm, in which at each iteration the server combines the current and past gradient information from each worker to compute next update, to solve convex problems with high dimensions."
- "'Notations. We use bold lower-case letters such as \mathbf{m} to represent vectors, lower-case letters such as m to represent scalars, and upper-case curlicue letters such as \mathcal{S} to represent sets. Aggregated vectors are denoted by a line over vectors such as $\overline{\mathbf{m}}$. Byzantine vectors are denoted by a tilde over vectors such as $\widetilde{\mathbf{m}}$. $\|\mathbf{m}\|$ denotes the Euclidean norm of \mathbf{m} . $|S|$ is the cardinality of set S . \odot denotes element-wise multiplication (Hadamard product). All operations between vectors are element-wise operations in this paper (except inner products of vectors)."
- o "
- o Idea to let nodes aggregate incoming gradient updates per node
 - [237]
 - Shows that we cannot detect Byzantine agents based on model parameter updates, but only based on gradient updates. However, this requires sending also the first moment/second moment parameters across the network every time => overhead
- New architecture
- Warm-start
 - o Note that this is different from federated transfer learning, see [144]
- Datasets used
- [238] => used for the PSI-CA



-
- Apart from the hash functions similar to [239] who either accidentally rediscovered SRA or forgot to add a proper reference
- [240]
 - “The server requires a minimum number of TCNs to be queried in order to prevent the user querying single TCNs and potentially identifying infected users. If users only have one received contactEventTCN, the list of encountered TCNs is padded with randomly generated TCNs, such that they can still query the server. This increases the chance of a false positive result. We introduce a limit on the query rate, such that the app can only send a limited number of requests to the server at a time in order to prevent brute force attacks. Then users have to wait a certain period of time until they can query the server anew. The server’s public key pk_S can additionally change for each combined hourly uploaded patient TCN set. This complicates a brute force attack further, because the attacker needs a different bloom filter of their encountered TCNs for each hourly dataset. Using this approach to private set intersection cardinality might reveal the number of encountered TCNs to the server. This is not the case if we allow the user to directly download all TCNs of infected people and check for matches on their phone. Using the latter method, the server gets no information whatsoever about non-infected app users (except maybe their IP-address), however the user could identify infected people using an attack as described

in Section 4.1. As an update to this protocol the bloom filter can be replaced with a cuckoo filter, which has the benefit of having lower error probabilities for the same size. This is also what is used in [241].”

- Commutative encryption algorithms, e.g. Pohlig-Hellman[242] or SRA[243]

Assumptions

- assumption: nodes hebben allemaal redelijk wat training data / geen backdoor attacks

Results

Illustrate performance under non-i.i.d. conditions

- 5 graphs with 2 classes per node, 4 classes per node, ... all classes per node

Illustrate power of direct averaging:

Map GARs to accuracy over time of first node, when the first node only generates very little data and other nodes generate much more data

Illustrate power of keeping track of last x updates:

Map GARs to accuracy over time of first node, when the other nodes only very seldom send an update

Illustrate power of exploration vs exploitation:

Map GARs to accuracy over time of first node, when it joins in quite late when the other nodes have already much better accuracy

Illustrate power of PSI-CA + shared dataset

Map GARs to accuracy over time of first node, when the data is utmost non-i.i.d. but still partly overlapping

Results when setting is asynchronous

Map GARs to accuracy over time () dataset 1 () 50 nodes () attack 1 () 75 attackers	Map GARs to accuracy over time () dataset 2 () 50 nodes () attack 1 () 75 attackers	Map GARs to accuracy over time () dataset 3 () 50 nodes () attack 1 () 75 attackers
Map GARs to accuracy over time () dataset 1 () 10 nodes () attack 1 () 75 attackers	Map GARs to accuracy over time () dataset 1 () 50 nodes () attack 1 () 75 attackers	Map GARs to accuracy over time () dataset 1 () 250 nodes () attack 1 () 75 attackers
Map GARs to accuracy over time () dataset 1 () 50 nodes () attack 1 () 75 attackers	Map GARs to accuracy over time () dataset 1 () 50 nodes () attack 2 () 75 attackers	Map GARs to accuracy over time () dataset 1 () 50 nodes () attack 3 () 75 attackers
Map GARs to accuracy over time () dataset 1 () 50 nodes () attack 1 () 10 attackers	Map GARs to accuracy over time () dataset 1 () 50 nodes () attack 1 () 75 attackers	Map GARs to accuracy over time () dataset 1 () 50 nodes () attack 1 () 175 attackers

Results when setting is non-i.i.d.

See table in previous section

Illustrate in mildly asynchronous / mildly non-i.i.d. setting

- Difference between simulated / distributed
- Impact of parameters (like exploration ratio)
- Impact of communication pattern

Discussion

Conclusion

References

1. Liu, Y., Ma, S., Aafer, Y., Lee, W.-C., Zhai, J., Wang, W., and Zhang, X., 'Trojaning Attack on Neural Networks', 2017.
2. Bagdasaryan, E., Veit, A., Hua, Y., Estrin, D., and Shmatikov, V., 'How to Backdoor Federated Learning', in, *International Conference on Artificial Intelligence and Statistics*, (PMLR, 2020)
3. Biggio, B., Didaci, L., Fumera, G., and Roli, F., 'Poisoning Attacks to Compromise Face Templates', in, *2013 International Conference on Biometrics (ICB)*, (IEEE, 2013)
4. Biggio, B., Nelson, B., and Laskov, P., 'Poisoning Attacks against Support Vector Machines', *arXiv preprint arXiv:1206.6389*, 2012.
5. Jagielski, M., Oprea, A., Biggio, B., Liu, C., Nita-Rotaru, C., and Li, B., 'Manipulating Machine Learning: Poisoning Attacks and Countermeasures for Regression Learning', in, *2018 IEEE Symposium on Security and Privacy (SP)*, (IEEE, 2018)
6. Li, B., Wang, Y., Singh, A., and Vorobeychik, Y., 'Data Poisoning Attacks on Factorization-Based Collaborative Filtering', in, *Advances in neural information processing systems*, (2016)
7. Rubinstein, B.I., Nelson, B., Huang, L., Joseph, A.D., Lau, S.-h., Rao, S., Taft, N., and Tygar, J.D., 'Antidote: Understanding and Defending against Poisoning of Anomaly Detectors', in, *Proceedings of the 9th ACM SIGCOMM conference on Internet measurement*, (2009)
8. Xiao, H., Biggio, B., Brown, G., Fumera, G., Eckert, C., and Roli, F., 'Is Feature Selection Secure against Training Data Poisoning?', in, *International Conference on Machine Learning*, (2015)
9. Yang, G., Gong, N.Z., and Cai, Y., 'Fake Co-Visitation Injection Attacks to Recommender Systems', in, *NDSS*, (2017)
10. Chen, X., Liu, C., Li, B., Lu, K., and Song, D., 'Targeted Backdoor Attacks on Deep Learning Systems Using Data Poisoning', *arXiv preprint arXiv:1712.05526*, 2017.
11. Koh, P.W. and Liang, P., 'Understanding Black-Box Predictions Via Influence Functions', *arXiv preprint arXiv:1703.04730*, 2017.
12. Suci, O., Marginean, R., Kaya, Y., Daume III, H., and Dumitras, T., 'When Does Machine Learning {Fail}? Generalized Transferability for Evasion and Poisoning Attacks', in, *27th {USENIX} Security Symposium ({USENIX} Security 18)*, (2018)
13. Bhagoji, A.N., Chakraborty, S., Mittal, P., and Calo, S., 'Analyzing Federated Learning through an Adversarial Lens', in, *International Conference on Machine Learning*, (PMLR, 2019)
14. Gu, T., Dolan-Gavitt, B., and Garg, S., 'Badnets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain', *arXiv preprint arXiv:1708.06733*, 2017.
15. Nelson, B., Barreno, M., Chi, F.J., Joseph, A.D., Rubinstein, B.I., Saini, U., Sutton, C.A., Tygar, J.D., and Xia, K., 'Exploiting Machine Learning to Subvert Your Spam Filter', *LEET*, 2008, 8, pp. 1-9.
16. Shafahi, A., Huang, W.R., Najibi, M., Suci, O., Studer, C., Dumitras, T., and Goldstein, T., 'Poison Frogs! Targeted Clean-Label Poisoning Attacks on Neural Networks', in, *Advances in neural information processing systems*, (2018)
17. Shen, S., Tople, S., and Saxena, P., 'Auror: Defending against Poisoning Attacks in Collaborative Deep Learning Systems', in, *Proceedings of the 32nd Annual Conference on Computer Security Applications*, (2016)
18. Baruch, G., Baruch, M., and Goldberg, Y., 'A Little Is Enough: Circumventing Defenses for Distributed Learning', in, *Advances in neural information processing systems*, (2019)
19. Bhagoji, A.N., Chakraborty, S., Mittal, P., and Calo, S., 'Model Poisoning Attacks in Federated Learning', in, *In Workshop on Security in Machine Learning (SecML), collocated with the 32nd Conference on Neural Information Processing Systems (NeurIPS'18)*, (2018)
20. Sun, Z., Kairouz, P., Suresh, A.T., and McMahan, H.B., 'Can You Really Backdoor Federated Learning?', *arXiv preprint arXiv:1911.07963*, 2019.
21. Xie, C., Huang, K., Chen, P.-Y., and Li, B., 'DbA: Distributed Backdoor Attacks against Federated Learning', in, *International Conference on Learning Representations*, (2019)
22. Zou, M., Shi, Y., Wang, C., Li, F., Song, W., and Wang, Y., 'Petrojan: Powerful Neural-Level Trojan Designs in Deep Learning Models', *arXiv preprint arXiv:1802.03043*, 2018.
23. Koloskova, A., Stich, S.U., and Jaggi, M., 'Decentralized Stochastic Optimization and Gossip Algorithms with Compressed Communication', *arXiv preprint arXiv:1902.00340*, 2019.
24. Lesort, T., Lomonaco, V., Stoian, A., Maltoni, D., Filliat, D., and Díaz-Rodríguez, N., 'Continual Learning for Robotics: Definition, Framework, Learning Strategies, Opportunities and Challenges', *Information fusion*, 2020, 58, pp. 52-68.
25. Verbraeken, J., Wolting, M., Katzy, J., Kloppenburg, J., Verbelen, T., and Rellermeyer, J.S., 'A Survey on Distributed Machine Learning', *ACM Computing Surveys (CSUR)*, 2020, 53, (2), pp. 1-33.
26. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., and Isard, M., 'Tensorflow: A System for Large-Scale Machine Learning', in, *12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16)*, (2016)
27. Medicare, C.f. and Medicaid Services, The Health Insurance Portability and Accountability Act of 1996 (Hippaa)', (1996)
28. Voigt, P. and Von dem Bussche, A., 'The Eu General Data Protection Regulation (Gdpr)', *A Practical Guide, 1st Ed.*, Cham: Springer International Publishing, 2017.
29. Chen, Z., Tian, P., Liao, W., and Yu, W., 'Zero Knowledge Clustering Based Adversarial Mitigation in Heterogeneous Federated Learning', *IEEE Transactions on Network Science and Engineering*, 2020.
30. McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B.A., 'Communication-Efficient Learning of Deep Networks from Decentralized Data', in, *Artificial Intelligence and Statistics*, (PMLR, 2017)
31. <https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>, accessed Date Accessed 2017 Accessed
32. Hard, A., Rao, K., Mathews, R., Ramaswamy, S., Beaufays, F., Augenstein, S., Eichner, H., Kiddon, C., and Ramage, D., 'Federated Learning for Mobile Keyboard Prediction', *arXiv preprint arXiv:1811.03604*, 2018.
33. Yang, T., Andrew, G., Eichner, H., Sun, H., Li, W., Kong, N., Ramage, D., and Beaufays, F., 'Applied Federated Learning:

- Improving Google Keyboard Query Suggestions', *arXiv preprint arXiv:1812.02903*, 2018.
34. Chen, M., Mathews, R., Ouyang, T., and Beaufays, F., 'Federated Learning of out-of-Vocabulary Words', *arXiv preprint arXiv:1903.10635*, 2019.
 35. Ramaswamy, S., Mathews, R., Rao, K., and Beaufays, F., 'Federated Learning for Emoji Prediction in a Mobile Keyboard', *arXiv preprint arXiv:1906.04329*, 2019.
 36. Chen, M., Suresh, A.T., Mathews, R., Wong, A., Allauzen, C., Beaufays, F., and Riley, M., 'Federated Learning of N-Gram Language Models', *arXiv preprint arXiv:1910.03432*, 2019.
 37. Yuan, B., Ge, S., and Xing, W., 'A Federated Learning Framework for Healthcare Iot Devices', *arXiv preprint arXiv:2005.05083*, 2020.
 38. Leroy, D., Coucke, A., Lavril, T., Gisselbrecht, T., and Dureau, J., 'Federated Learning for Keyword Spotting', in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, (IEEE, 2019)
 39. Sim, K.C., Beaufays, F., Benard, A., Guliani, D., Kabel, A., Khare, N., Lucassen, T., Zadrazil, P., Zhang, H., and Johnson, L., 'Personalization of End-to-End Speech Recognition on Mobile Devices for Named Entities', in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, (IEEE, 2019)
 40. Niknam, S., Dhillon, H.S., and Reed, J.H., 'Federated Learning for Wireless Communications: Motivation, Opportunities, and Challenges', *IEEE Communications Magazine*, 2020, 58, (6), pp. 46-51.
 41. Chen, M., Poor, H.V., Saad, W., and Cui, S., 'Wireless Communications for Collaborative Federated Learning in the Internet of Things', *arXiv preprint arXiv:2006.02499*, 2020.
 42. Lin, K.-Y. and Huang, W.-R., 'Using Federated Learning on Malware Classification', in *2020 22nd International Conference on Advanced Communication Technology (ICACT)*, (IEEE, 2020)
 43. Sozinov, K., Vlassov, V., and Girdzijauskas, S., 'Human Activity Recognition Using Federated Learning', in *2018 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Ubiquitous Computing & Communications, Big Data & Cloud Computing, Social Computing & Networking, Sustainable Computing & Communications (ISPA/IUCC/BDC/Cloud/SocialCom/SustainCom)*, (IEEE, 2018)
 44. Nguyen, T.D., Marchal, S., Miettinen, M., Fereidooni, H., Asokan, N., and Sadeghi, A.-R., 'Diot: A Federated Self-Learning Anomaly Detection System for Iot', in *2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS)*, (IEEE, 2019)
 45. Cetin, B., Lazar, A., Kim, J., Sim, A., and Wu, K., 'Federated Wireless Network Intrusion Detection', in *2019 IEEE International Conference on Big Data (Big Data)*, (IEEE, 2019)
 46. Lu, Y., Huang, X., Zhang, K., Maharjan, S., and Zhang, Y., 'Blockchain Empowered Asynchronous Federated Learning for Secure Data Sharing in Internet of Vehicles', *IEEE Transactions on Vehicular Technology*, 2020, 69, (4), pp. 4298-4311.
 47. Samarakoon, S., Bennis, M., Saad, W., and Debbah, M., 'Federated Learning for Ultra-Reliable Low-Latency V2v Communications', in *2018 IEEE Global Communications Conference (GLOBECOM)*, (IEEE, 2018)
 48. Gulati, A., Aujla, G.S., Chaudhary, R., Kumar, N., and Obaidat, M.S., 'Deep Learning-Based Content Centric Data Dissemination Scheme for Internet of Vehicles', in *2018 IEEE International Conference on Communications (ICC)*, (IEEE, 2018)
 49. Lu, Y., Huang, X., Dai, Y., Maharjan, S., and Zhang, Y., 'Federated Learning for Data Privacy Preservation in Vehicular Cyber-Physical Systems', *IEEE Network*, 2020, 34, (3), pp. 50-56.
 50. Liu, Y., James, J., Kang, J., Niyato, D., and Zhang, S., 'Privacy-Preserving Traffic Flow Prediction: A Federated Learning Approach', *IEEE Internet of Things Journal*, 2020.
 51. Mowla, N.I., Tran, N.H., Doh, I., and Chae, K., 'Federated Learning-Based Cognitive Detection of Jamming Attack in Flying Ad-Hoc Network', *IEEE Access*, 2019, 8, pp. 4338-4350.
 52. Liu, Y., Huang, A., Luo, Y., Huang, H., Liu, Y., Chen, Y., Feng, L., Chen, T., Yu, H., and Yang, Q., 'Fedvision: An Online Visual Object Detection Platform Powered by Federated Learning', in *AAAI*, (2020)
 53. Schneble, W. and Thamilarasu, G., 'Attack Detection Using Federated Learning in Medical Cyber-Physical Systems'.
 54. Lu, S., Zhang, Y., and Wang, Y., 'Decentralized Federated Learning for Electronic Health Records', in *2020 54th Annual Conference on Information Sciences and Systems (CISS)*, (IEEE, 2020)
 55. Brisimi, T.S., Chen, R., Mela, T., Olshevsky, A., Paschalidis, I.C., and Shi, W., 'Federated Learning of Predictive Models from Federated Electronic Health Records', *International journal of medical informatics*, 2018, 112, pp. 59-67.
 56. Xu, J. and Wang, F., 'Federated Learning for Healthcare Informatics', *arXiv preprint arXiv:1911.06270*, 2019.
 57. Rieke, N., Hancox, J., Li, W., Milletari, F., Roth, H., Albarqouni, S., Bakas, S., Galtier, M.N., Landman, B., and Maier-Hein, K., 'The Future of Digital Health with Federated Learning', *arXiv preprint arXiv:2003.08119*, 2020.
 58. Konečný, J., McMahan, H.B., Ramage, D., and Richtárik, P., 'Federated Optimization: Distributed Machine Learning for on-Device Intelligence', *arXiv preprint arXiv:1610.02527*, 2016.
 59. Bonawitz, K., Eichner, H., Grieskamp, W., Huba, D., Ingerman, A., Ivanov, V., Kiddon, C., Konečný, J., Mazzocchi, S., and McMahan, H.B., 'Towards Federated Learning at Scale: System Design', *arXiv preprint arXiv:1902.01046*, 2019.
 60. Lian, X., Zhang, C., Zhang, H., Hsieh, C.-J., Zhang, W., and Liu, J., 'Can Decentralized Algorithms Outperform Centralized Algorithms? A Case Study for Decentralized Parallel Stochastic Gradient Descent', in *Advances in neural information processing systems*, (2017)
 61. Xie, X., Ma, L., Wang, H., Li, Y., Liu, Y., and Li, X., 'Diffchaser: Detecting Disagreements for Deep Neural Networks', in *IJCAI*, (2019)
 62. Dobbe, R., Fridovich-Keil, D., and Tomlin, C., 'Fully Decentralized Policies for Multi-Agent Systems: An Information Theoretic Approach', in *Advances in neural information processing systems*, (2017)
 63. Tang, H., Gan, S., Zhang, C., Zhang, T., and Liu, J., 'Communication Compression for Decentralized Training', in *Advances in neural information processing systems*, (2018)
 64. Lalitha, A., Wang, X., Kilinc, O., Lu, Y., Javidi, T., and Koushanfar, F., 'Decentralized Bayesian Learning over Graphs', *arXiv preprint arXiv:1905.10466*, 2019.

65. Nedic, A. and Ozdaglar, A., 'Distributed Subgradient Methods for Multi-Agent Optimization', *IEEE Transactions on Automatic Control*, 2009, 54, (1), pp. 48-61.
66. Harinath, D., Satyanarayana, P., and Murthy, M., 'A Review on Security Issues and Attacks in Distributed Systems', *Journal of Advances in Information Technology*, 2017, 8, (1).
67. Lamport, L., Shostak, R., and Pease, M., 'The Byzantine Generals Problem', *Concurrency: The Works of Leslie Lamport*, (2019)
68. Blanchard, P., Guerraoui, R., and Stainer, J., 'Machine Learning with Adversaries: Byzantine Tolerant Gradient Descent', in, *Advances in neural information processing systems*, (2017)
69. Chen, Y., Su, L., and Xu, J., 'Distributed Statistical Machine Learning in Adversarial Settings: Byzantine Gradient Descent', *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 2017, 1, (2), pp. 1-25.
70. Zhang, Q., Cheng, L., and Boutaba, R., 'Algorithms and Architectures for Parallel Processing', *J. Int. Serv. Appl*, 2010, 1, (1), pp. 7-18.
71. El-Mhamdi, E.-M. and Guerraoui, R., 'Fast and Secure Distributed Learning in High Dimension', *arXiv preprint arXiv:1905.04374*, 2019.
72. Haykin, S., *Neural Networks and Learning Machines*, 3/E, (Pearson Education India, 2010)
73. Neelakantan, A., Vilnis, L., Le, Q.V., Sutskever, I., Kaiser, L., Kurach, K., and Martens, J., 'Adding Gradient Noise Improves Learning for Very Deep Networks', *arXiv preprint arXiv:1511.06807*, 2015.
74. Keskar, N.S., Mudigere, D., Nocedal, J., Smelyanskiy, M., and Tang, P.T.P., 'On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima', *arXiv preprint arXiv:1609.04836*, 2016.
75. Kleinberg, R., Li, Y., and Yuan, Y., 'An Alternative View: When Does Sgd Escape Local Minima?', *arXiv preprint arXiv:1802.06175*, 2018.
76. Bottou, L., 'Online Learning and Stochastic Approximations', *On-line learning in neural networks*, 1998, 17, (9), p. 142.
77. Chen, J., Pan, X., Monga, R., Bengio, S., and Jozefowicz, R., 'Revisiting Distributed Synchronous Sgd', *arXiv preprint arXiv:1604.00981*, 2016.
78. Wu, W., He, L., Lin, W., Mao, R., Maple, C., and Jarvis, S.A., 'Safa: A Semi-Asynchronous Protocol for Fast Federated Learning with Low Overhead', *IEEE Transactions on Computers*, 2020.
79. Xie, C., Koyejo, S., and Gupta, I., 'Asynchronous Federated Optimization', *arXiv preprint arXiv:1903.03934*, 2019.
80. Chen, Y., Ning, Y., and Rangwala, H., 'Asynchronous Online Federated Learning for Edge Devices', *arXiv preprint arXiv:1911.02134*, 2019.
81. Sprague, M.R., Jalalirad, A., Scavuzzo, M., Capota, C., Neun, M., Do, L., and Kopp, M., 'Asynchronous Federated Learning for Geospatial Applications', in, *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, (Springer, 2018)
82. Mohammad, U. and Sorour, S., 'Adaptive Task Allocation for Asynchronous Federated Mobile Edge Learning', *arXiv preprint arXiv:1905.01656*, 2019.
83. Yang, Y.-R. and Li, W.-J., 'Basgd: Buffered Asynchronous Sgd for Byzantine Learning', *arXiv preprint arXiv:2003.00937*, 2020.
84. Chen, M., Mao, B., and Ma, T., 'Efficient and Robust Asynchronous Federated Learning with Stragglers', in, *Submitted to International Conference on Learning Representations*, (2019)
85. Roy, A.G., Siddiqui, S., Pölsterl, S., Navab, N., and Wachinger, C., 'Braintorrent: A Peer-to-Peer Environment for Decentralized Federated Learning', *arXiv preprint arXiv:1905.06731*, 2019.
86. Hu, C., Jiang, J., and Wang, Z., 'Decentralized Federated Learning: A Segmented Gossip Approach', *arXiv preprint arXiv:1908.07782*, 2019.
87. Hegedűs, I., Danner, G., and Jelasity, M., 'Gossip Learning as a Decentralized Alternative to Federated Learning', in, *IFIP International Conference on Distributed Applications and Interoperable Systems*, (Springer, 2019)
88. Haseltalab, A. and Akar, M., 'Approximate Byzantine Consensus in Faulty Asynchronous Networks', in, *2015 American Control Conference (ACC)*, (IEEE, 2015)
89. Li, T., Sahu, A.K., Zaheer, M., Sanjabi, M., Talwalkar, A., and Smith, V., 'Federated Optimization in Heterogeneous Networks', *arXiv preprint arXiv:1812.06127*, 2018.
90. Nilsson, A., Smith, S., Ulm, G., Gustavsson, E., and Jirstrand, M., 'A Performance Evaluation of Federated Learning Algorithms', in, *Proceedings of the Second Workshop on Distributed Infrastructures for Deep Learning*, (2018)
91. Karimireddy, S.P., Kale, S., Mohri, M., Reddi, S.J., Stich, S.U., and Suresh, A.T., 'Scaffold: Stochastic Controlled Averaging for on-Device Federated Learning', *arXiv preprint arXiv:1910.06378*, 2019.
92. Robbins, H. and Monro, S., 'A Stochastic Approximation Method', *The annals of mathematical statistics*, 1951, pp. 400-407.
93. Kingma, D.P. and Ba, J., 'Adam: A Method for Stochastic Optimization', *arXiv preprint arXiv:1412.6980*, 2014.
94. Mukkamala, M.C. and Hein, M., 'Variants of Rmsprop and Adagrad with Logarithmic Regret Bounds', *arXiv preprint arXiv:1706.05507*, 2017.
95. Damaskinos, G., El Mhamdi, E.M., Guerraoui, R., Guirguis, A.H.A., and Rouault, S.L.A., 'Aggregathor: Byzantine Machine Learning Via Robust Gradient Aggregation', in, *The Conference on Systems and Machine Learning (SysML)*, 2019, (2019)
96. Zhang, S., Choromanska, A.E., and LeCun, Y., 'Deep Learning with Elastic Averaging Sgd', in, *Advances in neural information processing systems*, (2015)
97. Li, M., Andersen, D.G., Park, J.W., Smola, A.J., Ahmed, A., Josifovski, V., Long, J., Shekita, E.J., and Su, B.-Y., 'Scaling Distributed Machine Learning with the Parameter Server', in, *11th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 14)*, (2014)
98. Xing, E.P., Ho, Q., Xie, P., and Wei, D., 'Strategies and Principles of Distributed Machine Learning on Big Data', *Engineering*, 2016, 2, (2), pp. 179-195.
99. Yin, D., Chen, Y., Ramchandran, K., and Bartlett, P., 'Byzantine-Robust Distributed Learning: Towards Optimal Statistical Rates', *arXiv preprint arXiv:1803.01498*, 2018.

100. Muñoz-González, L., Biggio, B., Demontis, A., Paudice, A., Wongrassamee, V., Lupu, E.C., and Roli, F., 'Towards Poisoning of Deep Learning Algorithms with Back-Gradient Optimization', in, *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, (2017)
101. Wang, B. and Gong, N.Z., 'Attacking Graph-Based Classification Via Manipulating the Graph Structure', in, *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, (2019)
102. Mei, S. and Zhu, X., 'Using Machine Teaching to Identify Optimal Training-Set Attacks on Machine Learners', in, *AAAI*, (2015)
103. Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., Prakash, A., Kohno, T., and Song, D., 'Robust Physical-World Attacks on Deep Learning Visual Classification', in, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, (2018)
104. Fung, C., Yoon, C.J., and Beschastnikh, I., 'Mitigating Sybils in Federated Learning Poisoning', *arXiv preprint arXiv:1808.04866*, 2018.
105. Mhamdi, E.M.E., Guerraoui, R., and Rouault, S., 'The Hidden Vulnerability of Distributed Learning in Byzantium', *arXiv preprint arXiv:1802.07927*, 2018.
106. Fang, M., Cao, X., Jia, J., and Gong, N., 'Local Model Poisoning Attacks to Byzantine-Robust Federated Learning', in, *29th {USENIX} Security Symposium ({USENIX} Security 20)*, (2020)
107. Kairouz, P., McMahan, H.B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A.N., Bonawitz, K., Charles, Z., Cormode, G., and Cummings, R., 'Advances and Open Problems in Federated Learning', *arXiv preprint arXiv:1912.04977*, 2019.
108. Wang, H., Sreenivasan, K., Rajput, S., Vishwakarma, H., Agarwal, S., Sohn, J.-y., Lee, K., and Papailiopoulos, D., 'Attack of the Tails: Yes, You Really Can Backdoor Federated Learning', *arXiv preprint arXiv:2007.05084*, 2020.
109. Xie, C., Koyejo, O., and Gupta, I., 'Generalized Byzantine-Tolerant Sgd', *arXiv preprint arXiv:1802.10116*, 2018.
110. Guo, S., Zhang, T., Xie, X., Ma, L., Xiang, T., and Liu, Y., 'Towards Byzantine-Resilient Learning in Decentralized Systems', *arXiv preprint arXiv:2002.08569*, 2020.
111. Huber, P.J., *Robust Statistics*, (John Wiley & Sons, 2004)
112. Cretu, G.F., Stavrou, A., Locasto, M.E., Stolfo, S.J., and Keromytis, A.D., 'Casting out Demons: Sanitizing Training Data for Anomaly Sensors', in, *2008 IEEE Symposium on Security and Privacy (sp 2008)*, (IEEE, 2008)
113. Bhatia, K., Jain, P., and Kar, P., 'Robust Regression Via Hard Thresholding', in, *Advances in neural information processing systems*, (2015)
114. Diakonikolas, I., Kamath, G., Kane, D., Li, J., Moitra, A., and Stewart, A., 'Robust Estimators in High-Dimensions without the Computational Intractability', *SIAM Journal on Computing*, 2019, 48, (2), pp. 742-864.
115. Lai, K.A., Rao, A.B., and Vempala, S., 'Agnostic Estimation of Mean and Covariance', in, *2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS)*, (IEEE, 2016)
116. Yang, Z., Gang, A., and Bajwa, W.U., 'Adversary-Resilient Inference and Machine Learning: From Distributed to Decentralized', *stat*, 2019, 1050, p. 23.
117. Su, L. and Vaidya, N.H., 'Fault-Tolerant Distributed Optimization (Part Iv): Constrained Optimization with Arbitrary Directed Networks', *arXiv preprint arXiv:1511.01821*, 2015.
118. Sundaram, S. and Ghahserifard, B., 'Distributed Optimization under Adversarial Nodes', *IEEE Transactions on Automatic Control*, 2018, 64, (3), pp. 1063-1076.
119. Chen, L., Wang, H., Charles, Z., and Papailiopoulos, D., 'Draco: Byzantine-Resilient Distributed Training Via Redundant Gradients', *arXiv preprint arXiv:1803.09877*, 2018.
120. Alon, N., Matias, Y., and Szegedy, M., 'The Space Complexity of Approximating the Frequency Moments', *Journal of Computer and system sciences*, 1999, 58, (1), pp. 137-147.
121. Jerrum, M.R., Valiant, L.G., and Vazirani, V.V., 'Random Generation of Combinatorial Structures from a Uniform Distribution', *Theoretical computer science*, 1986, 43, pp. 169-188.
122. Lerasle, M. and Oliveira, R.I., 'Robust Empirical Mean Estimators', *arXiv preprint arXiv:1112.3914*, 2011.
123. Minsker, S., 'Geometric Median and Robust Estimation in Banach Spaces', *Bernoulli*, 2015, 21, (4), pp. 2308-2335.
124. Minsker, S., 'Distributed Statistical Estimation and Rates of Convergence in Normal Approximation', *Electronic Journal of Statistics*, 2019, 13, (2), pp. 5213-5252.
125. Bernstein, J., Wang, Y.-X., Azizzadenesheli, K., and Anandkumar, A., 'Signsgd: Compressed Optimisation for Non-Convex Problems', *arXiv preprint arXiv:1802.04434*, 2018.
126. Bernstein, J., Zhao, J., Azizzadenesheli, K., and Anandkumar, A., 'Signsgd with Majority Vote Is Communication Efficient and Fault Tolerant', *arXiv preprint arXiv:1810.05291*, 2018.
127. Chen, X., Chen, T., Sun, H., Wu, Z.S., and Hong, M., 'Distributed Training with Heterogeneous Data: Bridging Median- and Mean-Based Algorithms', *arXiv preprint arXiv:1906.01736*, 2019.
128. Sohn, J.-y., Han, D.-J., Choi, B., and Moon, J., 'Election Coding for Distributed Learning: Protecting Signsgd against Byzantine Attacks', *arXiv preprint arXiv:1910.06093*, 2019.
129. Li, L., Xu, W., Chen, T., Giannakis, G.B., and Ling, Q., 'Rsa: Byzantine-Robust Stochastic Aggregation Methods for Distributed Learning from Heterogeneous Datasets', in, *Proceedings of the AAAI Conference on Artificial Intelligence*, (2019)
130. Cao, D., Chang, S., Lin, Z., Liu, G., and Sun, D., 'Understanding Distributed Poisoning Attack in Federated Learning', in, *2019 IEEE 25th International Conference on Parallel and Distributed Systems (ICPADS)*, (IEEE, 2019)
131. Yang, Z. and Bajwa, W.U., 'Byrdie: Byzantine-Resilient Distributed Coordinate Descent for Decentralized Learning', *IEEE Transactions on Signal and Information Processing over Networks*, 2019, 5, (4), pp. 611-627.
132. Yang, Z. and Bajwa, W.U., 'Bridge: Byzantine-Resilient Decentralized Gradient Descent', *arXiv preprint arXiv:1908.08098*, 2019.
133. Peng, J. and Ling, Q., 'Byzantine-Robust Decentralized Stochastic Optimization', in, *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, (IEEE, 2020)
134. He, L., Karimireddy, S.P., and Jaggi, M., 'Byzantine-Robust Learning on Heterogeneous Datasets Via Resampling', *arXiv preprint arXiv:2006.09365*, 2020.

135. Gupta, N., Liu, S., and Vaidya, N.H., 'Byzantine Fault-Tolerant Distributed Machine Learning Using Stochastic Gradient Descent (Sgd) and Norm-Based Comparative Gradient Elimination (Cge)', *arXiv preprint arXiv:2008.04699*, 2020.
136. Barreno, M., Nelson, B., Joseph, A.D., and Tygar, J.D., 'The Security of Machine Learning', *Machine Learning*, 2010, 81, (2), pp. 121-148.
137. Tran, B., Li, J., and Madry, A., 'Spectral Signatures in Backdoor Attacks', in *Advances in neural information processing systems*, (2018)
138. Zhao, L., Hu, S., Wang, Q., Jiang, J., Chao, S., Luo, X., and Hu, P., 'Shielding Collaborative Learning: Mitigating Poisoning Attacks through Client-Side Detection', *IEEE Transactions on Dependable and Secure Computing*, 2020.
139. Jin, R., He, X., and Dai, H., 'Distributed Byzantine Tolerant Stochastic Gradient Descent in the Era of Big Data', in *ICC 2019-2019 IEEE International Conference on Communications (ICC)*, (IEEE, 2019)
140. Xie, C., Koyejo, S., and Gupta, I., 'Zeno: Distributed Stochastic Gradient Descent with Suspicion-Based Fault-Tolerance', in *International Conference on Machine Learning*, (PMLR, 2019)
141. Xie, C., Koyejo, S., and Gupta, I., 'Zeno++: Robust Fully Asynchronous Sgd', *arXiv preprint arXiv:1903.07020*, 2019.
142. Zhao, Y., Chen, J., Zhang, J., Wu, D., Teng, J., and Yu, S., 'Pdgan: A Novel Poisoning Defense Method in Federated Learning Using Generative Adversarial Network', in *International Conference on Algorithms and Architectures for Parallel Processing*, (Springer, 2019)
143. Wang, Z., Song, M., Zhang, Z., Song, Y., Wang, Q., and Qi, H., 'Beyond Inferring Class Representatives: User-Level Privacy Leakage from Federated Learning', in *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*, (IEEE, 2019)
144. Mothukuri, V., Parizi, R.M., Pouriyeh, S., Huang, Y., Dehghantanha, A., and Srivastava, G., 'A Survey on Security and Privacy of Federated Learning', *Future Generation Computer Systems*, 2020.
145. Liu, K., Dolan-Gavitt, B., and Garg, S., 'Fine-Pruning: Defending against Backdooring Attacks on Deep Neural Networks', in *International Symposium on Research in Attacks, Intrusions, and Defenses*, (Springer, 2018)
146. Wang, B., Yao, Y., Shan, S., Li, H., Viswanath, B., Zheng, H., and Zhao, B.Y., 'Neural Cleanse: Identifying and Mitigating Backdoor Attacks in Neural Networks', in *2019 IEEE Symposium on Security and Privacy (SP)*, (IEEE, 2019)
147. Jiang, Y., Wang, S., Ko, B.J., Lee, W.-H., and Tassioulas, L., 'Model Pruning Enables Efficient Federated Learning on Edge Devices', *arXiv preprint arXiv:1909.12326*, 2019.
148. Koh, P.W., Steinhardt, J., and Liang, P., 'Stronger Data Poisoning Attacks Break Data Sanitization Defenses', *arXiv preprint arXiv:1811.00741*, 2018.
149. Steinhardt, J., Koh, P.W., and Liang, P.S., 'Certified Defenses for Data Poisoning Attacks', in *Advances in neural information processing systems*, (2017)
150. Qiao, M. and Valiant, G., 'Learning Discrete Distributions from Untrusted Batches', *arXiv preprint arXiv:1711.08113*, 2017.
151. Chen, B., Carvalho, W., Baracaldo, N., Ludwig, H., Edwards, B., Lee, T., Molloy, I., and Srivastava, B., 'Detecting Backdoor Attacks on Deep Neural Networks by Activation Clustering', *arXiv preprint arXiv:1811.03728*, 2018.
152. Chou, E., Tramèr, F., Pellegrino, G., and Boneh, D., 'Sentinet: Detecting Physical Attacks against Deep Learning Systems', *arXiv preprint arXiv:1812.00292*, 2018.
153. Shen, Y. and Sanghavi, S., 'Learning with Bad Training Data Via Iterative Trimmed Loss Minimization', in *International Conference on Machine Learning*, (PMLR, 2019)
154. Diakonikolas, I., Kamath, G., Kane, D., Li, J., Steinhardt, J., and Stewart, A., 'Sever: A Robust Meta-Algorithm for Stochastic Optimization', in *International Conference on Machine Learning*, (2019)
155. Regatti, J. and Gupta, A., 'Befriending the Byzantines through Reputation Scores', *arXiv preprint arXiv:2006.13421*, 2020.
156. Azulay, S., Raz, L., Globerson, A., Koren, T., and Afek, Y., 'Holdout Sgd: Byzantine Tolerant Federated Learning', *arXiv preprint arXiv:2008.04612*, 2020.
157. Schmid, R., Pfitzner, B., Beilharz, J., Amrich, B., and Polze, A., 'Tangle Ledger for Decentralized Learning', in *2020 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, (IEEE, 2020)
158. Kim, H., Kim, S.-H., Hwang, J.Y., and Seo, C., 'Efficient Privacy-Preserving Machine Learning for Blockchain Network', *IEEE Access*, 2019, 7, pp. 136481-136495.
159. Shayan, M., Fung, C., Yoon, C.J., and Beschastnikh, I., 'Biscotti: A Ledger for Private and Secure Peer-to-Peer Machine Learning', *arXiv preprint arXiv:1811.09904*, 2018.
160. Chen, X., Ji, J., Luo, C., Liao, W., and Li, P., 'When Machine Learning Meets Blockchain: A Decentralized, Privacy-Preserving and Secure Design', in *2018 IEEE International Conference on Big Data (Big Data)*, (IEEE, 2018)
161. Kim, H., Park, J., Bennis, M., and Kim, S.-L., 'Blockchain on-Device Federated Learning', *IEEE Communications Letters*, 2019, 24, (6), pp. 1279-1283.
162. Kim, Y.J. and Hong, C.S., 'Blockchain-Based Node-Aware Dynamic Weighting Methods for Improving Federated Learning Performance', in *2019 20th Asia-Pacific Network Operations and Management Symposium (APNOMS)*, (IEEE, 2019)
163. Weng, J., Weng, J., Zhang, J., Li, M., Zhang, Y., and Luo, W., 'Deepchain: Auditable and Privacy-Preserving Deep Learning with Blockchain-Based Incentive', *IEEE Transactions on Dependable and Secure Computing*, 2019.
164. Zhou, S., Huang, H., Chen, W., Zhou, P., Zheng, Z., and Guo, S., 'Pirate: A Blockchain-Based Secure Framework of Distributed Machine Learning in 5g Networks', *IEEE Network*, 2020.
165. Toyoda, K. and Zhang, A.N., 'Mechanism Design for an Incentive-Aware Blockchain-Enabled Federated Learning Platform', in *2019 IEEE International Conference on Big Data (Big Data)*, (IEEE, 2019)
166. Majeed, U. and Hong, C.S., 'Flchain: Federated Learning Via Mec-Enabled Blockchain Network', in *2019 20th Asia-Pacific Network Operations and Management Symposium (APNOMS)*, (IEEE, 2019)
167. Salah, K., Rehman, M.H.U., Nizamuddin, N., and Al-Fuqaha, A., 'Blockchain for Ai: Review and Open Research Challenges', *IEEE Access*, 2019, 7, pp. 10127-10149.

168. Bao, X., Su, C., Xiong, Y., Huang, W., and Hu, Y., 'Flchain: A Blockchain for Auditable Federated Learning with Trust and Incentive', in *2019 5th International Conference on Big Data Computing and Communications (BIGCOM)*, (IEEE, 2019)
169. TOYODA, K., MATHIOPOULOS, P.T., and ZHANG, A.N., 'Novel Blockchain-Based Incentive-Aware Federated Learning Platform with Mechanism Design'.
170. Zhao, Y., Zhao, J., Jiang, L., Tan, R., and Niyato, D., 'Mobile Edge Computing, Blockchain and Reputation-Based Crowdsourcing Iot Federated Learning: A Secure, Decentralized and Privacy-Preserving System', *arXiv preprint arXiv:1906.10893*, 2019.
171. Kang, J., Xiong, Z., Niyato, D., Yu, H., Liang, Y.-C., and Kim, D.I., 'Incentive Design for Efficient Federated Learning in Mobile Networks: A Contract Theory Approach', in *2019 IEEE VTS Asia Pacific Wireless Communications Symposium (APWCS)*, (IEEE, 2019)
172. Kang, J., Xiong, Z., Niyato, D., Xie, S., and Zhang, J., 'Incentive Mechanism for Reliable Federated Learning: A Joint Optimization Approach to Combining Reputation and Contract Theory', *IEEE Internet of Things Journal*, 2019, 6, (6), pp. 10700-10714.
173. Zhao, Y., Zhao, J., Jiang, L., Tan, R., Niyato, D., Li, Z., Lyu, L., and Liu, Y., 'Privacy-Preserving Blockchain-Based Federated Learning for Iot Devices', *IEEE Internet of Things Journal*, 2020.
174. Preuveneers, D., Rimmer, V., Tsingenopoulos, I., Spooren, J., Joosen, W., and Ilie-Zudor, E., 'Chained Anomaly Detection Models for Federated Learning: An Intrusion Detection Case Study', *Applied Sciences*, 2018, 8, (12), p. 2663.
175. Gilad, Y., Hemo, R., Micali, S., Vlachos, G., and Zeldovich, N., 'Algorand: Scaling Byzantine Agreements for Cryptocurrencies', in *Proceedings of the 26th Symposium on Operating Systems Principles*, (2017)
176. Zhan, Y., Li, P., Qu, Z., Zeng, D., and Guo, S., 'A Learning-Based Incentive Mechanism for Federated Learning', *IEEE Internet of Things Journal*, 2020.
177. Khan, L.U., Tran, N.H., Pandey, S.R., Saad, W., Han, Z., Nguyen, M.N., and Hong, C.S., 'Federated Learning for Edge Networks: Resource Optimization and Incentive Mechanism', *arXiv preprint arXiv:1911.05642*, 2019.
178. Yu, H., Liu, Z., Liu, Y., Chen, T., Cong, M., Weng, X., Niyato, D., and Yang, Q., 'A Fairness-Aware Incentive Scheme for Federated Learning', in *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, (2020)
179. Hu, R. and Gong, Y., 'Trading Data for Learning: Incentive Mechanism for on-Device Federated Learning', *arXiv preprint arXiv:2009.05604*, 2020.
180. Zeng, R., Zhang, S., Wang, J., and Chu, X., 'Fmore: An Incentive Scheme of Multi-Dimensional Auction for Federated Learning in Mec', *arXiv preprint arXiv:2002.09699*, 2020.
181. Yu, H., Liu, Z., Liu, Y., Chen, T., Cong, M., Weng, X., Niyato, D., and Yang, Q., 'A Sustainable Incentive Scheme for Federated Learning', *IEEE Intelligent Systems*, 2020.
182. Le, T.H.T., Tran, N.H., Tun, Y.K., Nguyen, M.N., Pandey, S.R., Han, Z., and Hong, C.S., 'An Incentive Mechanism for Federated Learning in Wireless Cellular Network: An Auction Approach', *arXiv preprint arXiv:2009.10269*, 2020.
183. Cong, M., Yu, H., Weng, X., Qu, J., Liu, Y., and Yiu, S.M., 'A Vcg-Based Fair Incentive Mechanism for Federated Learning', *arXiv preprint arXiv:2008.06680*, 2020.
184. Ding, N., Fang, Z., and Huang, J., 'Incentive Mechanism Design for Federated Learning with Multi-Dimensional Private Information', in *2020 18th International Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOPT)*, (IEEE, 2020)
185. Lim, W.Y.B., Xiong, Z., Kang, J., Niyato, D., Zhang, Y., Leung, C., and Miao, C., 'An Incentive Scheme for Federated Learning in the Sky', in *Proceedings of the 2nd ACM MobiCom Workshop on Drone Assisted Wireless Communications for 5G and Beyond*, (2020)
186. Lim, W.Y.B., Xiong, Z., Miao, C., Niyato, D., Yang, Q., Leung, C., and Poor, H.V., 'Hierarchical Incentive Mechanism Design for Federated Machine Learning in Mobile Networks', *IEEE Internet of Things Journal*, 2020.
187. Ng, K.L., Chen, Z., Zelei Liu, H.Y., Liu, Y., and Yang, Q., 'A Multi-Player Game for Studying Federated Learning Incentive Schemes'.
188. Pandey, S.R., Suhail, S., Tun, Y.K., Alsenwi, M., and Hong, C.S., 'An Incentive Design to Perform Federated Learning'.
189. Feng, S., Niyato, D., Wang, P., Kim, D.I., and Liang, Y.-C., 'Joint Service Pricing and Cooperative Relay Communication for Federated Learning', in *2019 International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData)*, (IEEE, 2019)
190. Sarikaya, Y. and Ercetin, O., 'Motivating Workers in Federated Learning: A Stackelberg Game Perspective', *IEEE Networking Letters*, 2019, 2, (1), pp. 23-27.
191. Bonawitz, K., Ivanov, V., Kreuter, B., Marcedone, A., McMahan, H.B., Patel, S., Ramage, D., Segal, A., and Seth, K., 'Practical Secure Aggregation for Federated Learning on User-Held Data', *arXiv preprint arXiv:1611.04482*, 2016.
192. Sabt, M., Achemlal, M., and Bouabdallah, A., 'Trusted Execution Environment: What It Is, and What It Is Not', in *2015 IEEE Trustcom/BigDataSE/ISPA*, (IEEE, 2015)
193. Mo, F. and Haddadi, H., 'Efficient and Private Federated Learning Using Tee', in *EuroSys*, (2019)
194. Chen, Y., Luo, F., Li, T., Xiang, T., Liu, Z., and Li, J., 'A Training-Integrity Privacy-Preserving Federated Learning Scheme with Trusted Execution Environment', *Information Sciences*, 2020, 522, pp. 69-79.
195. Ji, J., Chen, X., Wang, Q., Yu, L., and Li, P., 'Learning to Learn Gradient Aggregation by Gradient Descent', in *IJCAI*, (2019)
196. Li, S., Cheng, Y., Wang, W., Liu, Y., and Chen, T., 'Learning to Detect Malicious Clients for Robust Federated Learning', *arXiv preprint arXiv:2002.00211*, 2020.
197. Rajput, S., Wang, H., Charles, Z., and Papailiopoulos, D., 'Detox: A Redundancy-Based Framework for Faster and More Robust Gradient Aggregation', in *Advances in neural information processing systems*, (2019)
198. Data, D., Song, L., and Diggavi, S., 'Data Encoding for Byzantine-Resilient Distributed Optimization', *arXiv preprint arXiv:1907.02664*, 2019.

199. Wang, S., Tuor, T., Salonidis, T., Leung, K.K., Makaya, C., He, T., and Chan, K., 'Adaptive Federated Learning in Resource Constrained Edge Computing Systems', *IEEE Journal on Selected Areas in Communications*, 2019, 37, (6), pp. 1205-1221.
200. Zhao, Y., Li, M., Lai, L., Suda, N., Civin, D., and Chandra, V., 'Federated Learning with Non-Iid Data', *arXiv preprint arXiv:1806.00582*, 2018.
201. Zhang, Y. and Yang, Q., 'A Survey on Multi-Task Learning', *arXiv preprint arXiv:1707.08114*, 2017.
202. Hertz, J.A., *Introduction to the Theory of Neural Computation*, (CRC Press, 2018)
203. Parisi, G.I., Tani, J., Weber, C., and Wermtner, S., 'Lifelong Learning of Human Actions with Deep Neural Network Self-Organization', *Neural Networks*, 2017, 96, pp. 137-149.
204. Parisi, G.I., Tani, J., Weber, C., and Wermtner, S., 'Lifelong Learning of Spatiotemporal Representations with Dual-Memory Recurrent Self-Organization', *Frontiers in neurobotics*, 2018, 12, p. 78.
205. Rabinowitz, N.C., Desjardins, G., Rusu, A.-A., Kavukcuoglu, K., Hadsell, R.T., Pascanu, R., Kirkpatrick, J., and Soyer, H.J., 'Progressive Neural Networks', (Google Patents, 2017)
206. Li, Z. and Hoiem, D., 'Learning without Forgetting', *IEEE transactions on pattern analysis and machine intelligence*, 2017, 40, (12), pp. 2935-2947.
207. Kemker, R., McClure, M., Abitino, A., Hayes, T., and Kanan, C., 'Measuring Catastrophic Forgetting in Neural Networks', in *Proceedings of the AAAI Conference on Artificial Intelligence*, (2018)
208. Liu, X., Masana, M., Herranz, L., Van de Weijer, J., Lopez, A.M., and Bagdanov, A.D., 'Rotate Your Networks: Better Weight Consolidation and Less Catastrophic Forgetting', in *2018 24th International Conference on Pattern Recognition (ICPR)*, (IEEE, 2018)
209. Ritter, H., Botev, A., and Barber, D., 'Online Structured Laplace Approximations for Overcoming Catastrophic Forgetting', *arXiv preprint arXiv:1805.07810*, 2018.
210. Lee, S.-W., Kim, J.-H., Jun, J., Ha, J.-W., and Zhang, B.-T., 'Overcoming Catastrophic Forgetting by Incremental Moment Matching', *arXiv preprint arXiv:1703.08475*, 2017.
211. Zenke, F., Poole, B., and Ganguli, S., 'Continual Learning through Synaptic Intelligence', in *International Conference on Machine Learning*, (PMLR, 2017)
212. Robins, A., 'Catastrophic Forgetting in Neural Networks: The Role of Rehearsal Mechanisms', in *Proceedings 1993 The First New Zealand International Two-Stream Conference on Artificial Neural Networks and Expert Systems*, (IEEE, 1993)
213. Robins, A., 'Catastrophic Forgetting, Rehearsal and Pseudorehearsal', *Connection Science*, 1995, 7, (2), pp. 123-146.
214. Gepperth, A. and Karaoguz, C., 'A Bio-Inspired Incremental Learning Architecture for Applied Perceptual Problems', *Cognitive Computation*, 2016, 8, (5), pp. 924-934.
215. Rebuffi, S.-A., Kolesnikov, A., Sperl, G., and Lampert, C.H., 'Icarl: Incremental Classifier and Representation Learning', in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, (2017)
216. Ma, C., Konečný, J., Jaggi, M., Smith, V., Jordan, M.I., Richtárik, P., and Takáč, M., 'Distributed Optimization with Arbitrary Local Solvers', *Optimization Methods and Software*, 2017, 32, (4), pp. 813-848.
217. Jaggi, M., Smith, V., Takáč, M., Terhorst, J., Krishnan, S., Hofmann, T., and Jordan, M.I., 'Communication-Efficient Distributed Dual Coordinate Ascent', *Advances in neural information processing systems*, 2014, 27, pp. 3068-3076.
218. Smith, V., Chiang, C.-K., Sanjabi, M., and Talwalkar, A.S., 'Federated Multi-Task Learning', in *Advances in neural information processing systems*, (2017)
219. Kumar, S., Dutta, S., Chaturvedi, S., and Bhatia, M., 'Strategies for Enhancing Training and Privacy in Blockchain Enabled Federated Learning', in *2020 IEEE Sixth International Conference on Multimedia Big Data (BigMM)*, (IEEE, 2020)
220. Kopparapu, K. and Lin, E., 'Fedfmc: Sequential Efficient Federated Learning on Non-Iid Data', *arXiv preprint arXiv:2006.10937*, 2020.
221. Yao, X. and Sun, L., 'Continual Local Training for Better Initialization of Federated Models', in *2020 IEEE International Conference on Image Processing (ICIP)*, (IEEE, 2020)
222. Gonzalez, C., Sakas, G., and Mukhopadhyay, A., 'What Is Wrong with Continual Learning in Medical Image Segmentation?', *arXiv preprint arXiv:2010.11008*, 2020.
223. Ling, C.X. and Bohn, T., 'A Conceptual Framework for Lifelong Learning', *arXiv preprint arXiv:1911.09704*, 2019.
224. Lomonaco, V., 'Continual Learning with Deep Architectures', 2019.
225. Shoham, N., Avidor, T., Keren, A., Israel, N., Benditkis, D., Mor-Yosef, L., and Zeitak, I., 'Overcoming Forgetting in Federated Learning on Non-Iid Data', *arXiv preprint arXiv:1910.07796*, 2019.
226. <https://www.kaggle.com/datasets?search=smartphone&sort=votes>, accessed Date Accessed
227. Stokkink, Q., Epema, D., and Pouwelse, J., 'A Truly Self-Sovereign Identity System', *arXiv preprint arXiv:2007.00415*, 2020.
228. Stokkink, Q. and Pouwelse, J., 'Deployment of a Blockchain-Based Self-Sovereign Identity', in *2018 IEEE international conference on Internet of Things (iThings) and IEEE green computing and communications (GreenCom) and IEEE cyber, physical and social computing (CPSCom) and IEEE smart data (SmartData)*, (IEEE, 2018)
229. Pouwelse, J.A., Garbacki, P., Wang, J., Bakker, A., Yang, J., Iosup, A., Epema, D.H., Reinders, M., Van Steen, M.R., and Sips, H.J., 'Tribler: A Social-Based Peer-to-Peer System', *Concurrency and computation: Practice and experience*, 2008, 20, (2), pp. 127-138.
230. Zeilemaker, N., Capotă, M., Bakker, A., and Pouwelse, J., 'Tribler: P2p Media Search and Sharing', in *Proceedings of the 19th ACM international conference on Multimedia*, (2011)
231. Boyd, S., Ghosh, A., Prabhakar, B., and Shah, D., 'Randomized Gossip Algorithms', *IEEE transactions on information theory*, 2006, 52, (6), pp. 2508-2530.
232. Jin, P.H., Yuan, Q., Iandola, F., and Keutzer, K., 'How to Scale Distributed Deep Learning?', *arXiv preprint arXiv:1611.04581*, 2016.
233. Chang, H., Shejwalkar, V., Shokri, R., and Houmansadr, A., 'Cronus: Robust and Heterogeneous Collaborative Learning with

- Black-Box Knowledge Transfer', *arXiv preprint arXiv:1912.11279*, 2019.
234. Xie, C., Koyejo, O., and Gupta, I., 'Fall of Empires: Breaking Byzantine-Tolerant Sgd by Inner Product Manipulation', in, *Uncertainty in Artificial Intelligence*, (PMLR, 2020)
235. Cao, X. and Lai, L., 'Distributed Gradient Descent Algorithm Robust to an Arbitrary Number of Byzantine Attackers', *IEEE Transactions on Signal Processing*, 2019, 67, (22), pp. 5850-5864.
236. Alistarh, D., Allen-Zhu, Z., and Li, J., 'Byzantine Stochastic Gradient Descent', in, *Advances in neural information processing systems*, (2018)
237. Chen, C., Zhang, J., Tung, A.K., Kankanhalli, M., and Chen, G., 'Robust Federated Recommendation System', *arXiv preprint arXiv:2006.08259*, 2020.
238. Lv, S., Ye, J., Yin, S., Cheng, X., Feng, C., Liu, X., Li, R., Li, Z., Liu, Z., and Zhou, L., 'Unbalanced Private Set Intersection Cardinality Protocol with Low Communication Cost', *Future Generation Computer Systems*, 2020, 102, pp. 1054-1061.
239. De Cristofaro, E., Gasti, P., and Tsudik, G., 'Fast and Private Computation of Cardinality of Set Intersection and Union', in, *International Conference on Cryptology and Network Security*, (Springer, 2012)
240. Holzappel, K., Karl, M., Lotz, L., Carle, G., Djefal, C., Fruck, C., Haack, C., Heckmann, D., Kindt, P.H., and Köppl, M., 'Digital Contact Tracing Service: An Improved Decentralized Design for Privacy and Effectiveness', *arXiv preprint arXiv:2006.16960*, 2020.
241. Kales, D., Rechberger, C., Schneider, T., Senker, M., and Weinert, C., 'Mobile Private Contact Discovery at Scale', in, *28th {USENIX} Security Symposium ({USENIX} Security 19)*, (2019)
242. Pohlig, S. and Hellman, M., 'An Improved Algorithm for Computing Logarithms over $G_f(P)$ and Its Cryptographic Significance (Corresp.)', *IEEE transactions on information theory*, 1978, 24, (1), pp. 106-110.
243. Shamir, A., Rivest, R.L., and Adleman, L.M., 'Mental Poker', *The Mathematical Gardner*, (Springer, 1981)