# Delta Sum Learning: an approach for fast and global convergence in Gossip Learning

Tom Goethals, Merlijn Sebrechts, Stijn De Schrijver, Filip De Turck and Bruno Volckaert

*Ghent University - imec, IDLab*

Gent, Belgium

ORCID: 0000-0002-1332-2290, 0000-0002-4093-7338, N/A, 0000-0003-4824-1199, 0000-0003-0575-5894

*Abstract*—Federated Learning is a popular approach for distributed learning due to its security and computational benefits. With the advent of powerful devices in the network edge, Gossip Learning further decentralizes Federated Learning by removing centralized integration and relying fully on peer to peer updates. However, the averaging methods generally used in both Federated and Gossip Learning are not ideal for model accuracy and global convergence. Additionally, there are few options to deploy Learning workloads in the edge as part of a larger application using a declarative approach such as Kubernetes manifests. This paper proposes Delta Sum Learning as a method to improve the basic aggregation operation in Gossip Learning, and implements it in a decentralized orchestration framework based on Open Application Model, which allows for dynamic node discovery and intent-driven deployment of multi-workload applications. Evaluation results show that Delta Sum performance is on par with alternative integration methods for 10 node topologies, but results in a 58% lower global accuracy drop when scaling to 50 nodes. Overall, it shows strong global convergence and a logarithmic loss of accuracy with increasing topology size compared to a linear loss for alternatives under limited connectivity.

*Index Terms*—gossip learning, artificial intelligence, ai, edge learning, federated learning

## I. INTRODUCTION

Federated Learning (FL) has since long improved the training of Artificial Intelligence (AI) models by enabling distributed training on huge datasets while integrating the results at a central location, resulting in reduced training times and practical use of larger datasets. This is achieved through any number of integration methods, most popularly averaging methods (e.g. Federated Averaging or FedAvg). With the rise of edge computing, there is a push to further decentralize AI training through Gossip Learning (GL), which eliminates the centralized aggregation and relies on each node to spread its updates to the rest of the cluster by proxy, thus gossip. GL may be useful for various reasons; data may be processed at the source to reduce privacy concerns, spare computational capacity from dedicated hardware may be leveraged, etc. However, GL has certain inefficiencies compared to FL, for example, updates must be sent to several nodes instead of a single centralized location. Furthermore, because the distribution of training data and received updates is likely asymmetrical,

models are likely to diverge during training. As a result, GL requires the integration of the full model whereas some FL approaches only send changed parameters to reduce network traffic. To counteract this, GL nodes should send updates to as few others as possible, which in turn leads to slower or possibly no global convergence, especially if not all nodes in a gossip cluster are mutually known. Furthermore, standard averaging approaches may result in statistical anomalies (e.g. vanishing variance), which lead to slower training convergence, as well as lower accuracy. As such, there is a need for integration methods that provide strong global convergence, through minimal and local communication with other nodes.

In terms of framework support, ML workloads are strongly supported in cloud environments. For example, KubeFlow enables MLOps in Kubernetes clusters, automating workload deployment and management using Kubernetes manifests on up to hundreds of nodes. However, edge learning lacks the same support; solutions such as Azure IoT Edge and OpenEI generally only enable edge intelligence, and in limited cases edge learning, whereas most recent studies focus on edge learning specifically through FL, or are aimed at specialized use cases that only leverage learning workloads. As such, there is a clear need for a framework that enables GL in the edge as a part of larger applications, while allowing Kubernetes-like modeling of workloads.

This paper examines the fundamental integration operation used by FL and GL, proposing Delta Sum Learning to improve global convergence. Additionally, this method is integrated into a decentralized intent-based framework for edge workload deployment, including Gossip and ML services to support GL.

Concretely, the contributions of this paper are:

- Improving Gossip Learning through Delta Sum Learning to provide better and faster global convergence with minimal local communication.
- Integrating the improvements into a prototype framework for online edge learning and distributed updates.
- Evaluating the efficiency and resource impact of the improvements compared to other Gossip Learning approaches.

The rest of this article is organized as follows: Section II presents existing research related to the topic, while Section III introduces a mathematical model to improve Gossip Learning, which is integrated into a practical framework in Section IV. Section V details the evaluation setup and scenarios, while

the results are presented in Section VI. Section VII presents topics for future work, and finally, Section VIII draws high level conclusions from the paper.

## II. RELATED WORK

Federated Learning (FL) and by extension Decentralized Federated Learning (DFL) have been extensively studied [1] for their potential in Neural Network training. Specifically, a variety of model averaging or aggregation methods can be used, such as Federated Averaging, various Secure Aggregation [2] methods which mitigate privacy risks by obfuscating individual model updates, or Matched Averaging [3] methods aiming to reduce communication overhead. Furthermore, the effect of network topologies [4] and data distribution [5] on the performance of (D)FL has been extensively studied, as well as mitigating noisy communications channels [6].

Furthermore, a number of methods have been proposed that improve specific security aspects of (D)FL, including Byzantine-Robust FL [7] for adversarial node behavior, or Distributionally Robust averaging [8] which considers Non-IID data properties. Additionally, some methods aim to improve specific integration aspects, such as vanishing variance issues [9]. Finally, Xiao et al. argue that straightforward averaging methods are far from ideal for distributed learning [10], showing an increasing correlation between nodes while computed distances between their models stop diminishing.

Gossip Learning [11] (GL) has been proposed as a fully decentralized variant of DFL, using gossip algorithms to spread model updates throughout complex topologies. In GL, each node only directly contacts a small set of neighbours to push updates, usually through a gossip protocol such as SWIM [12]. Gossip protocols may also benefit from push-pull operations [13], which is shown to offer superior update dissemination compared to push only as in most DFL approaches. Inheriting many (dis)advantages from DFL, its performance is shown to be on par with FL depending on hyperparameters, network data compression and integration methods [14].

Network topology can have a significant impact on both (D)FL and GL; apart from the topologies discussed in [4], more efficient update schemes can be constructed for FL when nodes are grouped in different silos [15], and the decentralized nature of GL is shown to be robust in dynamic network environments [16].

Existing frameworks such as EdgeFL [17] offer a variety of options to run FL algorithms in the network edge, however, unlike the framework proposed in this paper they do not present a generic orchestration framework which allows ML workloads and GL to be used as components of a larger application. There are a number of frameworks that solve more specific FL issues. For example, DeFTA considers the network outdegree of contributing nodes in DFL to weigh their updates [18], and integrates it into a DFL framework. Fedstellar [19] offers a similar framework, with various options to run DFL workloads, showing its effectiveness in attack detection.

## III. DELTA SUM LEARNING

This section builds a model for the improvement of the weight update integration step in Gossip Learning (GL), starting from the fundamental mechanism of batch updates on a single and multiple nodes.

A single batch of training on one node can be represented as:

$$w_{t+1} = w_t + \alpha \sum_{i=0}^{S} \frac{\partial F(x_i, w_t, y_i)}{\partial w_t} \tag{1}$$

Where $w_t$ are the model weights at time $t$, $\alpha$ is the learning rate, $S$ iterates over the samples in a batch, and $F(x_i, w_t, y_i)$ is the model loss function using sample input $x_i$ and expected output (e.g. classification) $y_i$. Explaining this intuitively, the partial derivative is the slope vector at the hypersurface point defined by the model loss function $F(...)$ as an outcome of the model weights $w_t$, with respect to the weights $w_t$, which is exactly the factor required for gradient descent. Although multiple model layers are not explicitly represented here, $F(...)$ can represent any and all layers.

Extending this to multiple training epochs $T$, as well as a group of $N$ nodes simultaneously training on parts of the same dataset (i.e. FL), the expected integration becomes:

$$w_{t_0+T} = w_{t_0} + \alpha \sum_{t=t_0}^{t_0+T} \sum_{n=0}^{N} \sum_{i=0}^{S} \frac{\partial F(x_{n,i}, w_{n,t}, y_{n,i})}{\partial w_t} \tag{2}$$

While there is a discrepancy here in the slightly different weights matrices $w_{n,t}$ learned by each node compared to $w_t$ in monolithic learning, for FL there is little difference between this and training in (large) batches, as the weight matrix of each node can be periodically synchronized with the centralized main model to keep $w_{n,t}$ from diverging. This is not the case with GL, however there is an issue with averaging methods to be considered first. In averaging approaches using full model updates, an update $u_{r,t}$ is received from each node every round, and integrated into a main model. Although FL often works by sending only weight tensor deltas instead of the entire model, the nature of GL requires the entire model be sent as there is no global model to integrate into; for now full model updates are assumed. The total received update $u_{a,t}$ is then:

$$u_{a,t} = u_{l,t} + u_{r,t} \tag{3}$$

Where $u_{l,t}$ is the local learning update (if any), and $u_{r,t}$ is the combined update from remote nodes. These can be expressed in terms of Eq. 1, resulting in (respectively):

$$u_{l,t} = w_{a,t} + \alpha \sum_{i=0}^{S} \frac{\partial F(x_i, w_t, y_i)}{\partial w_t} \tag{4}$$

$$u_{r,t} = \sum_{n=0}^{N} \left( w_{n,t} + \alpha \sum_{i=0}^{S} \frac{\partial F(x_{n,i}, w_{n,t}, y_{n,i})}{\partial w_{n,t}} \right) \tag{5}$$

Or, when updates are the result of multiple learning epochs $T$, the analog of Eq. 2 becomes:

$$u_{a,t_0+T} = w_{a,t_0} + \alpha \sum_{t=t_0}^{t_0+T} \sum_{i=0}^{S} \frac{\partial F(x_i, w_{a,t}, y_i)}{\partial w_{a,t}}$$
$$+ \sum_{n=0}^{N} \left( w_{n,t_0} + \alpha \sum_{t=t_0}^{t_0+T} \sum_{i=0}^{S} \frac{\partial F(x_{n,i}, w_{n,t}, y_{n,i})}{\partial w_{n,t}} \right) \tag{6}$$

Using this for an update of node 0 from epoch 0 to $T$, since all models $w_{n,t_0}$ are the same as $w_{a,t_0}$ at $t = 0$, this reduces to:

$$u_{0,T} = w_{0,0}(N+1) + \alpha \sum_{t=0}^{T} \sum_{i=0}^{S} \frac{\partial F(x_i, w_{0,t}, y_i)}{\partial w_{0,t}}$$
$$+ \alpha \sum_{n=0}^{N} \sum_{t=0}^{T} \sum_{i=0}^{S} \frac{\partial F(x_{n,i}, w_{0,t}, y_{n,i})}{\partial w_{0,t}} \tag{7}$$

Or, substituting local and remote learning update factors with $\delta_{l,T}$ and $\delta_{r,T}$, respectively:

$$u_{0,T} = w_{0,0}(N+1) + (\delta_{l,T} + \delta_{r,T}) \tag{8}$$

Which, when applying any averaging method results in a new model weight tensor $w_{0,T}$:

$$w_{0,T} = w_{0,0} + \frac{(\delta_{l,T} + \delta_{r,T})}{N+1} \tag{9}$$

However, this shows that any straightforward averaging method imposes an implicit learning penalty based on the number of node updates received, as the actual updates are implicitly divided by $N + 1$. Although this approach assumes updates with full model weight tensors, the result is analogous for FL with federated averaging using only model weight updates:

$$w_{t+1} = w_t + \frac{\sum_{n=0}^{N} k_n \Delta w_{n,t}}{\sum_{n=0}^{N} k_n} \tag{10}$$

This method allows for different numbers of samples $k_n$ per node per update, while expressing the sum over multiple epochs $T$ from Eqs. 1 and 2 as $\Delta w_{n,t}$ for readability. Expressed differently:

$$w_{t+1} = w_t + \frac{\sum_{n=0}^{N} k_n \Delta w_{n,t}}{\mu_{kn} N} \tag{11}$$

Where $\mu_{kn}$ is the average number of samples per node. Eq. 2 implicitly has the same number of learning samples for every batch and node; to allow for different numbers of learning samples per node it can be rewritten in similar terms as Eq. 11:

$$w_{t+1} = w_t + \frac{\sum_{n=0}^{N} k_n \Delta w_{n,t}}{\mu_{kn}} \tag{12}$$

Or, extending $\Delta w_{n,t}$ back to its full form:

$$w_{t_0+T} = w_{t_0} + \frac{\sum_{n=0}^{N} k_n \sum_{t=t_0}^{t_0+T} \sum_{i=0}^{S} \alpha \frac{\partial F(x_{n,i}, w_{n,t}, y_{n,i})}{\partial w_t}}{\mu_{kn}} \tag{13}$$

This approach still weighs node updates relative to the training samples for each update, while avoiding $N$ in the denominator. However, due to differences in local model weights during training in GL, full model updates must be considered for further progress. To eliminate the $N + 1$ factor from the updates themselves in Eq. 9 and achieve a similar outcome to Eq. 2, the base weight tensor and update tensor for each node must be separately transmitted during training updates, allowing for:

$$w_{0,T} = w_{0,0} + (\delta_{l,T} + \delta_{r,T}) \tag{14}$$

Or, expanding the $\delta$ factors again and expressed in a more general form:

$$w_{a,t_0+T} = \frac{w_{a,t_0} + \sum_{n=0}^{N} w_{n,t_0}}{N+1} + \alpha \sum_{t=t_0}^{t_0+T} \sum_{i=0}^{S} \frac{\partial F(x_i, w_{a,t}, y_i)}{\partial w_{a,t}}$$
$$+ \alpha \sum_{n=0}^{N} \sum_{t=t_0}^{t_0+T} \sum_{i=0}^{S} \frac{\partial F(x_{n,i}, w_{n,t}, y_{n,i})}{\partial w_{n,t}} \tag{15}$$

As the interiors of both the local (double) and remote (triple) sums represent a number of epochs on the same nodes n, they can be summarized as $w_{n,t_0+T} - w_{n,t_0}$ (with $\alpha$ already integrated in the learning algorithm):

$$w_{a,t_0+T} = \frac{w_{a,t_0} + \sum_{n=0}^{N} w_{n,t_0}}{N+1} + (w_{a,t_0+T} - w_{a,t_0})$$
$$+ \sum_{n=0}^{N} (w_{n,t_0+T} - w_{n,t_0}) \tag{16}$$

Or, when interpreting the local weight matrix and delta as node "$N + 1$":

$$w_{a,t_0+T} = \frac{\sum_{n=0}^{N+1} w_{n,t_0}}{N+1} + \sum_{n=0}^{N+1} \Delta w_{n,t_0+T} \tag{17}$$

Eq. 6 shows that the models diverge after the first gossip update, especially since not all nodes are in direct contact. As a result, models will train into slightly different directions depending on their own data subset and initial updates received, especially during the first training rounds when accuracy and loss change rapidly when starting from randomized weight tensors. Therefore, all $w_{n,t}$ are likely to be minutely different throughout the learning process, although they should converge as updates spread through the gossip network, and they are assumed compatible due to training for the same global objective. Despite this, the derived model in Eq. 17 shows the

necessity to integrate the base weight matrices from all nodes for which updates were received through averaging; Eq. 2 uses only the local matrix because all nodes are assumed to train on the same base matrix (periodically updated if necessary).

Finally, there are some practical concerns. Because data and learning are not ideally spread in GL, there may be significant accuracy and loss jitter when receiving updates from nodes with contradictory updates. To reduce these effects, it is prudent to introduce a gossip learning factor $\lambda$:

$$w_{a,t_0+T} = \frac{\sum_{n=0}^{N+1} w_{n,t_0}}{N+1} + \lambda \sum_{n=0}^{N+1} \Delta w_{n,t_0+T} \qquad (18)$$

This factor should be small enough to dampen erratic updates, yet large enough that its effect does not fall below $1/(N+1)$, which would slow down learning more than standard averaging methods already do. An acceptable value is determined in the evaluation section. However, as learning topologies grow, nodes possess a relatively decreasing subset of the total learning set to train on; if training from local data does not perfectly align with the global optimum then training will have an increasingly detrimental effect on gossip updates, potentially canceling them out as:

$$w_{a,t_0} + \Delta w_{a,t_0+T} = \frac{\sum_{n=0}^{N} w_{n,t_0}}{N} + \sum_{n=0}^{N} (\Delta w_{n,t_0+T}) \quad (19)$$

Which may be partially avoided using a dynamic learning factor; models are assumed to stabilize throughout the topology over time, and the combined "wisdom" of the cluster should be increasingly prioritized over the deleterious effect of local learning. As such, the method of Delta Sum Learning is summarized as:

$$w_{a,t_0+T} = \frac{\sum_{n=0}^{N+1} w_{n,t_0}}{N+1} + \lambda(t_0+T) \sum_{n=0}^{N+1} \Delta w_{n,t_0+T} \quad (20)$$

$$\lambda(t) = min\left(A + \frac{t}{B}, C\right) \qquad (21)$$

Where A, B and C are constants to be empirically determined for optimal training. While B is strongly dependent on the expected maximum number of training epochs up to which the factor should increase, which may not be calculable during continuous or online training, functional values for both A and C are determined for the scenario in the evaluation section.

## IV. Architecture

To enable gossip cluster building and to evaluate Delta Sum Learning in a realistic scenario, the model is implemented in the architecture illustrated by Fig. 1. This architecture uses Flocky [20] as a basic framework as its services allow for the discovery of other nodes, as well as providing a decentralized cluster-wide metadata store and orchestration (Swirly and deployment) services.

Flocky is an open source framework that leverages Open Application Model (OAM) to model both node capabilities and workload requirements. Each node monitors its own resources, capabilities, and running workloads (i.e. Traits and Components, respectively) through Capability Providers (e.g. Workload agent and ML Service in Fig. 1), which are stored in the Repository Service and shared with nodes ① found by the Discovery Service. Not all nodes in a cluster are mutually known; the Discovery Service is configured to explore nodes only within certain parameters (e.g. latency, physical distance, ...), so each node only directly knows a small part of the whole cluster. Workloads are requested through the Swirly Service, which matches requirements with known node capabilities from the Repository Service, and subsequently deploys the workloads on the most suitable node through its Deployment Service, which in turn leverages a workload agent (e.g. Docker[1] or Feather [21]) to run them.

To enable Gossip Learning within Flocky, this research adds two additional services to manage model updates through gossiping:

The **ML Service** monitors running ML workloads for model updates consisting of a base model and a delta tensor ②. Currently, this is implemented through shared memory to improve performance, although a REST endpoint may be less error prone at the cost of possibly significant overhead (i.e. model serialization). Additionally, this service acts as a Capability Provider, mapping the ML capabilities of the local node and reporting them to the Repository Service. When a model update is detected, it is reported to the Gossip Service along with the workload identifier for further dissemination in the cluster. Conversely, whenever the Gossip Service reports a model update for a specific workload, it is retrieved and passed along with the source node identifier to the workload for integration. This service is currently implemented in Golang; for performance reasons, integration methods are implemented in the workloads themselves.

The **Gossip Service** leverages a Golang implementation of the SWIM gossip protocol to enable gossiping ③. Whereas the Discovery and Repository services constitute a gossip-like cluster, these are only used for OAM-based node and workload data, while the Gossip Service may be used for generic and operational workload data. Furthermore, using SWIM guarantees the propagation of workload data to relevant nodes, whereas the base Flocky services are designed to be error tolerant and "best effort". Unlike normal gossip clusters, the Gossip Service cluster is not created through configuration or joining; the nodes discovered by the Discovery Service are periodically synchronized as comprising the gossip cluster. The Gossip Service is generic: when a message is received (i.e. model update from ML Service), it forwards the message using its key (i.e. workload identifier) to other nodes in the cluster. Listeners can subscribe to the Gossip Service by providing a list of keys, e.g. the ML Service providing a list of its monitored workloads.

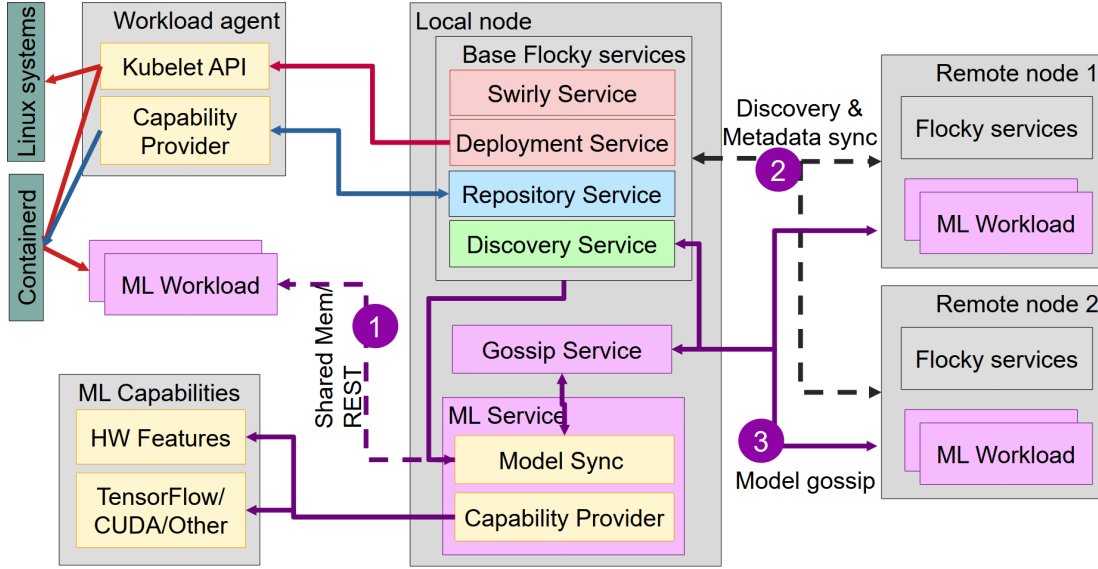---

[1]Docker Containers - https://www.docker.com/

Fig. 1: Architecture overview of decentralized cluster discovery and Gossip Learning showing: (1) Node discovery and metadata synchronization between nodes. (2) ML model synchronization between workloads and a dedicated service. (3) Model gossiping based on locally discovered nodes and hosted ML workloads.

The use of these services is enabled by adding a specific Trait to the Flocky framework, allowing workloads to be designated as Gossip Learning ML workloads. This Trait requires an endpoint within the workload (shared memory file or REST service) to exchange model updates, and when the Deployment Service applies this Trait it notifies the ML service of a new workload to be monitored.

## V. EVALUATION

This section describes the evaluation setup and any implementation details that deviate from Section IV, as well as the evaluation scenario and methodology. The code for all services and custom evaluation tools is made available on GitHub[2].

### A. Evaluation Setup

Evaluations are performed on a single Ubuntu 24.04 desktop PC with an Intel Core Ultra 7 155H processor, 32GiB dual channel LPDDR5, and a 1TB M.2 NVMe drive.

### B. Evaluation Scenario

The effectiveness of Delta Sum Learning is measured by comparing it to standard model averaging and variance corrected [9] averaging. Although many other integration methods exist, Delta Sum Learning can theoretically enhance every method that can be modified for separate base model and update dissemination. For performance reasons, all alternatives are implemented in the Python ML workload itself. Semi-random topologies are generated containing 10, 25 and 50 nodes, inspected to avoid bisection and to ensure that each node has between 1 and 8 neighbours for direct communication (Discovery and gossip). Average connectivity for the

topologies is determined to be 3.3, 3.2 and 4.2 neighbours, respectively. Each node runs Discovery, Repository, and Gossip/ML services, and an ML workload. The ML workload consists of a Convolutional Neural Network (CNN) comprised of two iterations of convolution and max pooling, followed by two dense layers, and uses the default MNIST digits classification dataset [22] for training. The model contains 55,338 parameters, for a total size of 221,352 bytes at fp32. The ML workload runs in training mode for 200 epochs, delivering model updates to the ML service, with integrations every 20 epochs. After training mode, it continues in convergence mode until round 235 to gauge the effect of additional integration rounds. During convergence rounds, full models are averaged rather than using Delta Sum Learning (or other strategies), as there is no longer any model delta.

### C. Methodology

All nodes in an evaluation scenario are simulated on the same machine by assigning incremental service ports, allowing discovery at localhost rather than distinct IP addresses. Training evaluation metrics are exported by the ML workload after each training epoch, and for each gossip update round after epoch 200. While absolute simultaneity is not guaranteed by the evaluation setup, the results show that equal CPU prioritization and timing individual processes using sleep operations is enough to ensure limited divergence in epoch count at any time across all simulated nodes. Each simulated node is assigned an equal share of the MNIST dataset, divided into a training portion and validation portion, as well as the full validation dataset to gauge global convergence. The dynamic gossip learning factor $\lambda$ constructed in Eq. 21 is used during training; for the MNIST scenario, the constants $A = 0.15$, $B = 1000$, and $C = 0.35$ are used. Both lower ($A = 0.05$

and C = 0.1) and higher (A = 0.35 and C = 0.6) learning factors result in worse performance for the specified scenario, although metaparameter optimization may provide slightly different values. Memory use was determined to be significantly higher when using hardware acceleration ($\sim$1.5GiB) than when running models on CPU only ($\sim$900MiB). Although some memory is shared across processes, the ML workloads are run entirely on CPU to ensure the feasibility of 50 node topologies. No artificial delays are used other than a short sleep operation intended to synchronize epochs with rounds; other than that, training epochs are executed as fast as possible, making evaluation strongly CPU bound. As such, CPU use is not monitored, while measured network traffic in larger topologies is not reliable due to significant throttling of gossip traffic resulting from CPU starvation. As a result, gossip rounds tend to lag training epochs in larger topologies. For all evaluations, CPU scaling is disabled to ensure a proper comparison. For practical purposes and to avoid timing issues, the ML workloads are run as host processes rather than containers. As a result, the Workload agent is stubbed, while the rest of the components (e.g. Discovery and Repository) function normally.

## VI. RESULTS

This section presents the results of the evaluations, considering both training metrics and networking aspects. Although network throughput was not explicitly measured due to issues discussed in Section V-C, theoretical performance can be compared to FL. For training, only the accuracy metric is represented. Although other metrics such as loss and F1 scores may provide a better indication for most models, loss was observed to be completely (inversely) in line with accuracy[3], and accuracy is sufficient to gauge a limited model.

### A. Accuracy and Convergence

Fig. 2 shows the accuracy of Delta Sum Learning over the entire runtime of the evaluation scenarios compared to standard averaging and variance corrected averaging. For small 10-node topologies, the difference between all strategies is negleg-ible; variance corrected averaging has a higher performance during the early stages but all models end up with identical accuracy a few integration rounds after the last training epoch. Importantly, accuracy ends up at 99.1%, providing a baseline for larger topologies. At 25 nodes, Delta Sum Learning results in a higher accuracy overall. Integrations every 20 epochs have a higher positive impact, while the negative effect of local training is less pronounced due to the dynamic gossip learning factor. The latter advantage is shared by variance corrected averaging, although standard averaging and variance corrected averaging result in equal accuracy at round 235. The 50 node topology illustrates the advantages of Delta Sum Learning in terms of scaling and global convergence. Despite each node having a minute share of the full training dataset, at epoch 50 Delta Training already results in a significantly

[3]https://github.com/idlab-discover/flocky/evaldata

higher accuracy than the alternatives, with integrations having a higher impact than in the 25 node topology. Additionally, local training has a significantly lower negative impact on global validation dataset accuracy. Furthermore, with Delta Sum Learning it only takes a few convergence rounds after training for median accuracy in the topology to improve to the top performing nodes, while the worst performing nodes catch up by round 230. While other strategies also converge to some degree, they retain a significant range above and below median accuracy beyond the end of the evaluation.

Fig. 3 compares the results at round 235 across topology sizes. Both standard averaging and variance corrected averaging sustain a linear accuracy loss as topologies grow, with variance corrected averaging showing a more uniform accuracy across larger topologies. Delta Sum Learning, however, has the same accuracy as the alternative strategies in a 10 node topology, but loses accuracy logarithmically as topology size grows: absolute accuracy drops from 99.1% at 10 nodes to 98.6% at 50 nodes, while the alternatives end up at 97.9%. Importantly, the results also show that the effect of the dynamic learning factor from Eq. 21 is independent of node connectivity; while the 25 node topology has a lower average connectivity and the the 50 node topology a higher connectivity than the 10 node topology, Delta Sum Learning performs increasingly better than alternatives with a growing number of nodes.

### B. Network Efficiency

Various expected trends for network efficiency are shown in Fig. 4, expressed in model updates per second as a model-agnostic unit. Due to the discovery algorithm, network traffic is expected to scale as $O(r^2\delta_n)$, with discovery distance r and local node density $\delta_n$, although the gossip protocol may be configured to limit the number of updates sent per round. Calculated trends are based on the Measured values for the 10 node topology; as illustrated by the chart the 25 node topology is significantly CPU limited, resulting in a real-time throughput drop as CPU intensive training takes longer to complete. Due to the manual checks imposed on the evaluation topologies, connectivity is relatively consistent across all scenarios, resulting in the Expected throughput adhering closely to expected throughput for Constant connectivity (i.e. either a physically growing topology, or $r^2 = 1/\delta_n$ if density increases). The trend for Connectivity increase indicates throughput for increasing node density in a topology of constant physical size and $r$. Finally, the expected network load for FedAvg is shown based on a model update every 5 seconds, and a model synchronization every 20 updates, in line with the 10 node GL scenario.

## VII. FUTURE WORK

The results show that Delta Sum Learning results in sigifni-cantly higher accuracy and better global convergence than standard averaging or variance corrected averaging applied to Gossip Learning. Furthermore, Delta Sum Learning exhibits a far lower accuracy loss than alternatives as topologies grow,
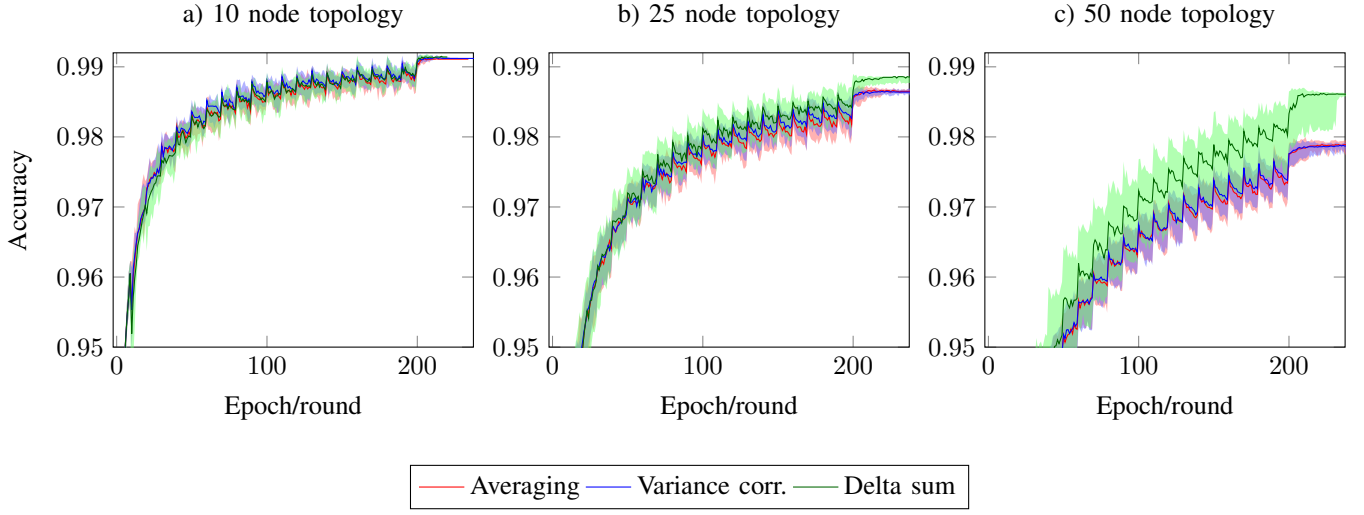
Fig. 2: Accuracy for Delta Sum Learning and alternative strategies for topology sizes ranging from 10 (a)) to 50 (c)). Shaded areas indicate the full range of accuracies of all nodes, except statistical outliers for visualization purposes. Comparing the different topology sizes, Delta Sum Learning results in a far lower accuracy loss as gossip topology size increases than alternative strategies.
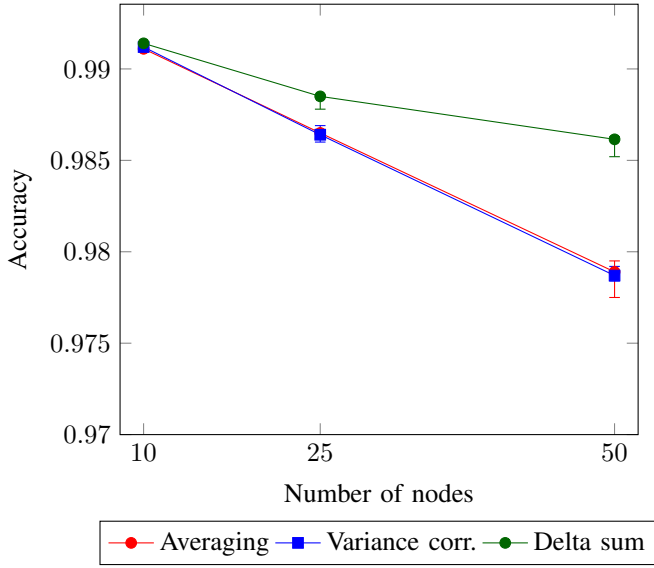


Fig. 3: Median accuracy of Delta Sum Learning compared to other strategies at round 235, for an increasing number of nodes in the topology. While accuracy loss for other strategies is almost linear with number of nodes, Delta Sum appears to follow a logarithmic pattern.
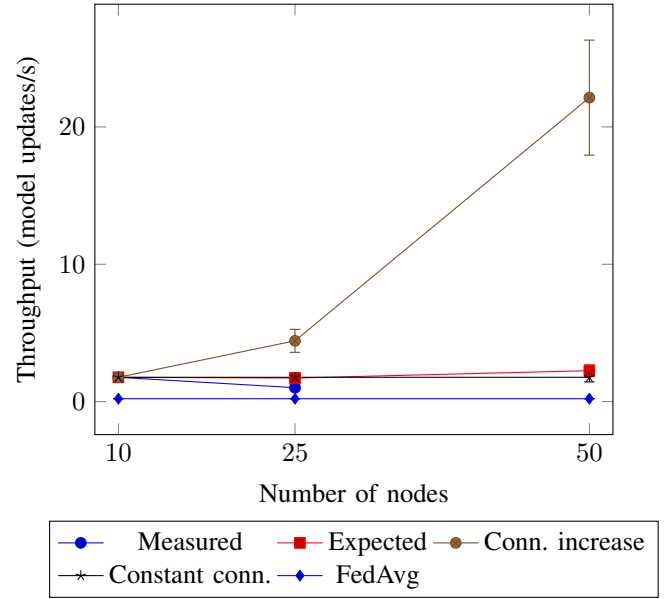


Fig. 4: Network throughput in model updates per second as measured, and calculated for various theoretical scenarios such as Constant connectivity, Connectivity increase due to node density, and FedAvg under the same conditions as the evaluations.

despite limited connectivity and model updates. However, some aspects of the proposed solution may be further improved.

The Delta Sum Learning model examines the impact of model integration from first hop (i.e. direct neighbour) updates. However, it does not take into account training data asymmetry and model divergence beyond direct neighbours; instead, it uses a gossip learning factor to mitigate these issues. While

this factor is also necessary to regulate the balance between local learning and remote updates, further examination of the impact of connectivity, model divergence during training, and data asymmetry may improve the model.

GL avoids bottlenecks by removing centralized integration, technically allowing scaling to larger training topologies than FL. Conversely, the evaluation shows GL requires significantly

more network traffic than FL in all (theorized) scenarios. While decentralized systems that require global consistency necessarily have a significant communication overhead, the absolute overhead may be reduced by, for example, modifying the gossip protocol to limit the number of updates sent per training epoch depending on connectivity, or limiting updates based on when model integrations are scheduled to occur. Furthermore, the Flocky discovery mechanism can be configured for lower connectivity, and the gossip protocol itself may be amenable to improvements for the specific use case of GL.

Finally, the gossip service may be improved to create a separate cluster per model, or per namespace (i.e. group of models or workload metadata). This approach avoids synchronization of specific data with nodes that have no use for it, which improves network efficiency and reduces security risks by restricting data to only those nodes that require it. In terms of security the implemented method of Delta Learning borrows some concepts from Secure Aggregation, as updates obfuscate any training data and, as with reducing network overhead, could be averaged over several epochs instead of distributing each epoch update to further reduce low-level information leakage.

## VIII. CONCLUSION

This paper introduces Delta Sum Learning to improve global convergence and training accuracy. It provides a mathematical model to characterize the operation and benefits of Delta Sum Learning, and presents a practical implementation which operates in Flocky, a decentralized orchestration framework. The evaluation section compares the implementation to standard averaging and variance-corrected averaging GL approaches using a basic CNN trained on the MNIST digits classification dataset. Evolution results indicate that Delta Sum Learning performance is increasingly better than the alternative approaches as topologies grow, resulting in a logarithmic loss of accuracy rather than a linear loss, as well as providing global convergence close to the best performing nodes instead of converging to median cases. Concretely, the global accuracy loss for Delta Sum Learning is 58% lower than the alternatives when scaling from 10 node to 50 node topologies. However, results also show that GL requires around 5 times more traffic between nodes to achieve global synchronization than FL, which scales with node connectivity but may be reduced through future work. Finally, some topics for future work are discussed concerning the accuracy, network overhead and security aspects of Delta Sum Learning and its implementation in Flocky.

## REFERENCES

[1] L. Yuan, Z. Wang, L. Sun, P. S. Yu, and C. G. Brinton, "Decentralized federated learning: A survey and perspective," *IEEE Internet of Things Journal*, vol. 11, no. 21, pp. 34 617–34 638, 2024.

[2] H. Fereidooni, S. Marchal, M. Miettinen, A. Mirhoseini, H. Möllering, T. D. Nguyen, P. Rieger, A.-R. Sadeghi, T. Schneider, H. Yalame, and S. Zeitouni, "Safelearn: Secure aggregation for private federated learning," in *2021 IEEE Security and Privacy Workshops (SPW)*, 2021, pp. 56–62.

[3] S. Shukla and N. Srivastava, "Federated matched averaging with information-gain based parameter sampling," in *Proceedings of the First International Conference on AI-ML Systems*, ser. AIMLSystems '21. New York, NY, USA: Association for Computing Machinery, 2021. [Online]. Available: https://doi.org/10.1145/3486001.3486225

[4] S. Wang, D. Li, J. Geng, Y. Gu, and Y. Cheng, "Impact of network topology on the performance of dml: Theoretical analysis and practical factors," in *IEEE INFOCOM 2019-IEEE conference on computer communications*. IEEE, 2019, pp. 1729–1737.

[5] H. Zhu, J. Xu, S. Liu, and Y. Jin, "Federated learning on non-iid data: A survey," *Neurocomputing*, vol. 465, pp. 371–390, 2021.

[6] V. P. Chellapandi, A. Upadhyay, A. Hashemi, and S. H. Żak, "Decentralized federated learning: Model update tracking under imperfect information sharing," in *2024 IEEE International Conference on Big Data (BigData)*, 2024, pp. 7697–7706.

[7] M. Fang, Z. Zhang, Hairi, P. Khanduri, J. Liu, S. Lu, Y. Liu, and N. Gong, "Byzantine-robust decentralized federated learning," in *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '24. New York, NY, USA: Association for Computing Machinery, 2024, p. 2874–2888. [Online]. Available: https://doi.org/10.1145/3658644.3670307

[8] Y. Deng, M. M. Kamani, and M. Mahdavi, "Distributionally robust federated averaging," in *Proceedings of the 34th International Conference on Neural Information Processing Systems*, ser. NIPS '20. Red Hook, NY, USA: Curran Associates Inc., 2020.

[9] Y. Tian, Z. Al-Ars, M. Kitsak, and P. Hofstee, "Vanishing variance problem in fully decentralized neural-network systems," 2024. [Online]. Available: https://arxiv.org/abs/2404.04616

[10] P. Xiao, S. Cheng, V. Stankovic, and D. Vukobratovic, "Averaging is probably not the optimum way of aggregating parameters in federated learning," *Entropy*, vol. 22, no. 3, 2020. [Online]. Available: https://www.mdpi.com/1099-4300/22/3/314

[11] L. Giaretta and S. Girdzijauskas, "Gossip learning: Off the beaten path," in *2019 IEEE International Conference on Big Data (Big Data)*, 2019, pp. 1117–1124.

[12] A. Das, I. Gupta, and A. Motivala, "Swim: scalable weakly-consistent infection-style process group membership protocol," in *Proceedings International Conference on Dependable Systems and Networks*, 2002, pp. 303–312.

[13] A. Srivastava, T. J. Maranzatto, and S. Ulukus, "Age of gossip with the push-pull protocol," in *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2025, pp. 1–5.

[14] I. Hegedűs, G. Danner, and M. Jelasity, "Decentralized learning works: An empirical comparison of gossip learning and federated learning," *Journal of Parallel and Distributed Computing*, vol. 148, pp. 109–124, 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0743731520303890

[15] O. MARFOQ, C. XU, G. Neglia, and R. Vidal, "Throughput-optimal topology design for cross-silo federated learning," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 19 478–19 487.

[16] A. Di Maio, M. A. Dinani, and G. Rizzo, "The upsides of turbulence: Baselining gossip learning in dynamic settings," in *Proceedings of the Twenty-fourth International Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing*, 2023, pp. 376–381.

[17] H. Zhang, J. Bosch, and H. H. Olsson, "Enabling efficient and low-effort decentralized federated learning with the edgefl framework," *Information and Software Technology*, vol. 178, p. 107600, 2025. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0950584924002052

[18] Y. Zhou, M. Shi, Y. Tian, Q. Ye, and J. Lv, "Defta: A plug-and-play peer-to-peer decentralized federated learning framework," *Information Sciences*, vol. 670, p. 120582, 2024. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S002002552400495X

[19] E. T. Martínez Beltrán, Ángel Luis Perales Gómez, C. Feng, P. M. Sánchez Sánchez, S. López Bernal, G. Bovet, M. Gil Pérez, G. Martínez Pérez, and A. Huertas Celdrán, "Fedstellar: A platform for decentralized federated learning," *Expert Systems with Applications*, vol. 242, p. 122861, 2024. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0957417423033638

[20] T. Goethals, M. Sebrechts, Mays Al-Naday, F. De Turck, and B. Volckaert, "Flocky: Decentralized intent-based edge

orchestration using open application model," 2025. [Online]. Available: 10.13140/RG.2.2.35590.48967

[21] T. Goethals, M. De Clercq, M. Sebrechts, F. De Turck, and B. Volckaert, "Feather: Lightweight container alternatives for deploying workloads in the edge." in *CLOSER*, 2024, pp. 27–37.

[22] U. N. I. of Standards and Technology, "Mnist database of handwritten digits," 1998.