

SwarmAvg : A Novel Approach to Fully Distributed Machine Learning

Josh Pattman

April 18, 2023



Abstract

TODO: This needs to be slightly reworded and also wordcounted

Federated learning is a technique that allows a machine learning model to be trained on data distributed across multiple data islands. This approach protects privacy by keeping the data decentralized, meaning that sensitive data does not need to ever be shared. Swarm learning is a similar technique that eliminates the need for a central server, ensuring that not only the data, but also the communication, is completely decentralized. In this paper, a novel swarm learning technique called SwarmAvg is presented which operates without a blockchain. The algorithm is validated against federated learning in various scenarios, and also in some situations where federated learning cannot be applied.

Statement of Originality

- I have read and understood the [ECS Academic Integrity](#) information and the University's [Academic Integrity Guidance for Students](#).
- I am aware that failure to act in accordance with the [Regulations Governing Academic Integrity](#) may lead to the imposition of penalties which, for the most serious cases, may include termination of programme.
- I consent to the University copying and distributing any or all of my work in any form and using third parties (who may be based outside the EU/EEA) to verify whether my work contains plagiarised material, and for quality assurance purposes.

You must change the statements in the boxes if you do not agree with them.

We expect you to acknowledge all sources of information (e.g. ideas, algorithms, data) using citations. You must also put quotation marks around any sections of text that you have copied without paraphrasing. If any figures or tables have been taken or modified from another source, you must explain this in the caption and cite the original source.

I have acknowledged all sources, and identified any content taken from elsewhere.

If you have used any code (e.g. open-source code), reference designs, or similar resources that have been produced by anyone else, you must list them in the box below. In the report, you must explain what was used and how it relates to the work you have done.

I have not used any resources produced by anyone else.

You can consult with module teaching staff/demonstrators, but you should not show anyone else your work (this includes uploading your work to publicly-accessible repositories e.g. Github, unless expressly permitted by the module leader), or help them to do theirs. For individual assignments, we expect you to work on your own. For group assignments, we expect that you work only with your allocated group. You must get permission in writing from the module teaching staff before you seek outside assistance, e.g. a proofreading service, and declare it here.

I did all the work myself, or with my allocated group, and have not helped anyone else.

We expect that you have not fabricated, modified or distorted any data, evidence, references, experimental results, or other material used or presented in the report. You must clearly describe your experiments and how the results were obtained, and include all data, source code and/or designs (either in the report, or submitted as a separate file) so that your results could be reproduced.

The material in the report is genuine, and I have included all my data/code/designs.

We expect that you have not previously submitted any part of this work for another assessment. You must get permission in writing from the module teaching staff before re-using any of your previously submitted work for this assessment.

I have not submitted any part of this work for another assessment.

If your work involved research/studies (including surveys) on human participants, their cells or data, or on animals, you must have been granted ethical approval before the work was carried out, and any experiments must have followed these requirements. You must give details of this in the report, and list the ethical approval reference number(s) in the box below.

My work did not involve human participants, their cells or data, or animals.

Contents

1	Background	6
1.1	Federated Learning	7
1.2	Swarm Learning	7
1.3	Blockchain	9
2	Problem and Goals	10
2.1	Problem	11
2.1.1	Privacy	11
2.1.2	Performance	11
2.2	Goals	11
2.2.1	Design an Novel Algorithm for Swarm Learning	11
2.2.2	Implement the Algorithm	12
2.2.3	Test Performance in Situations Where FedAvg Per- forms Well	12
2.2.4	Test Performance in Situations Where FedAvg Per- forms Badly	12
3	Analysis and Specification of the Solution to the Problem	13
4	Detailed Design	14
4.1	Node	15
4.2	Disposal of Blockchain	15
4.3	Training Counter	16
4.4	Model Combination Methods	16
4.4.1	Averaging	16
4.4.2	Averaging With Synchronisation Rate	17
4.4.3	Filtering By Training Counter	17
4.5	Sparse Network Behaviour	18
4.5.1	Passive Convergence	18
4.5.2	Relay	18
4.6	Complete Algorithm	19

5	Implementation	20
5.1	Dataset and Machine Learning Model	21
5.1.1	Dataset	21
5.1.2	Model	21
5.2	Federated Learning	21
5.2.1	Algorithm	22
5.2.2	Evaluation	22
5.3	Prototype	22
5.3.1	Algorithm	23
5.3.2	Evaluation	23
5.4	Final	23
5.4.1	Algorithm	23
5.4.2	Evaluation	24
6	Testing Strategy and Results	25
6.1	Methods	26
6.1.1	Metrics	26
6.1.2	Data Collection	26
6.1.3	Node Counts	26
6.1.4	Algorithm Configurations	26
6.1.5	Data Volume Per Node	26
6.1.6	Epochs	27
6.2	Dense Network Performance	27
6.2.1	Results	28
6.2.2	Analysis	30
6.3	Dense Network Performance with Node Dropout	31
6.4	Sparse Network Performance	31
6.4.1	Results	33
6.4.2	Analysis	35
6.5	SwarmAvg on a FedAvg topology	36
7	Critical Evaluation	37
7.1	Comparison of SL to FL	38
7.1.1	Advantages of using SL over FL	38
7.1.2	Disadvantages of using SL over FL	38
7.1.3	Limitations of the Presented Implementation	38
7.1.4	Further Extensions of Presented Implementation	38
8	Project Management	39
8.1	Time Management	40
8.2	Changing Plans	40

8.3	Risk Assessment	41
8.3.1	Personal Issues	41
8.3.2	Hardware Failure	41
8.3.3	Algorithm Does Not Work Work	42
9	Conclusion and Future Work	43
9.1	Conclusion	44
9.2	Future Work	44
A	Sward Learning Algorithm	45
A.1	Training Step	45
A.2	Model Received Event	46
B	Machine Learning Model	47
C	Gantt Charts	48
C.1	Gantt - Interim	49
C.2	Gantt - Final	50

Chapter 1

Background

DRAFT

1.1 Federated Learning

Many modern machine learning algorithms require large volumes of diverse data to achieve optimal performance. Real-world data is often distributed among multiple nodes that are unable to share it with each other due to privacy regulations such as GDPR [1], reducing the amount of data available to train models on and negatively impacting trained performance [2]. Federated Learning (FL) [3] is a technique in machine learning that aims to train a single model using all available data across nodes without requiring any data to be shared among them.

There are a multitude of published FL frameworks [4], each with different merits and drawbacks for certain use cases. Federated Averaging (FedAvg) [5] is a commonly used yet simple framework, which splits training into iterations where three steps take place:

1. A copy of the current model is sent to each node from the central server.
2. Each node performs some training with their copy of the model and their own private data.
3. The trained models from each node are sent back to the server to aggregate into the new server model.

The server model is improved over time, beyond what could be achieved by simply training on a single nodes data.

As it does not require data to be shared between nodes, FL is naturally beneficial for privacy sensitive tasks compared to conventional machine learning where the data is aggregated in a central location **CITE THIS STATEMENT**. Additionally, as FL performs training on multiple nodes in parallel, it can make better use of available training resources in situations where processing power is distributed among multiple nodes **CITE THIS STATEMENT**.

One branch of FL is distributed federated learning (DFL). This algorithm works in a very similar fashion to FL, but replaces the central server with an elected leader in a network of nodes [6]. DFL has been shown to increase fault tolerance and security over FL.

1.2 Swarm Learning

Swarm learning (SL) is a subcategory of FL which operates in a completely distributed and decentralised manner. SL enables the collaboration of nodes

to learn a shared global model, however in contrast to FL, a central server is never used. SL also does not use leader election, so all nodes on the network are given the same importance.

In SL, the model on which new training is performed is known as the global model. However, unlike FL where the global model is stored in a central location, the global model in SL does not materially exist, but is instead a concept which is agreed upon by the nodes in the network.

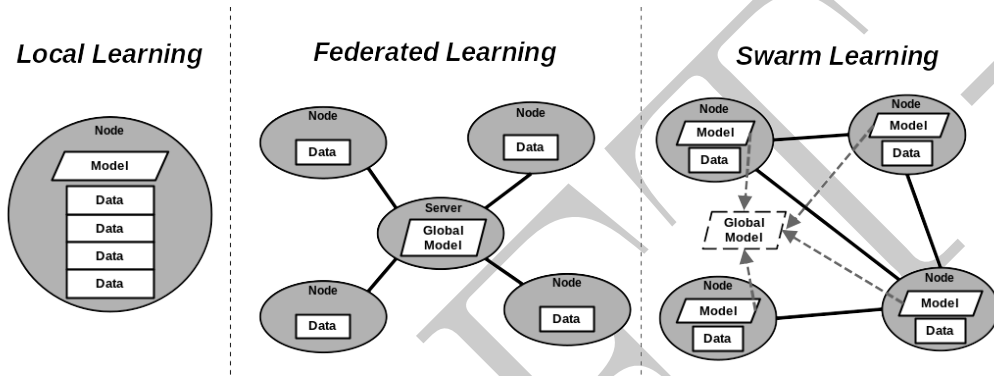


Figure 1.1: Diagram of different learning algorithms. Each *Node* indicates a single training machine, and each line denotes a connection between two machines, along which the model can be shared. In the swarm learning diagram, the dashed lines show that each local model is an approximation of the global model.

One SL algorithm, simply referred to as *Swarm Learning* by its creators, uses a blockchain to store the global model [7]. In this version of SL, training is performed by repeating the following steps:

1. A node obtains a copy of the global model from the blockchain.
2. The node performs some training with their copy of the model and their own private data.
3. The updated model is merged with the latest blockchain model and sent back to the blockchain for other nodes to use.

SL exhibits many of the same benefits of FL over conventional learning, but it also improves upon FL in some aspects. As is often the case when comparing decentralised algorithms to their centralised counterparts [8], the absence of a central server theoretically makes SL more resilient to failures than centralized FL approaches such as FedAvg. The lack of a leader election protocol also means that SL may be better suited to tasks where networks

of nodes are sparsely connected, as leader election in dynamic networks is a complex problem and can add large amounts of overhead [9]. Additionally, the removal of the need for a server in SL reduces the likelihood of performance bottlenecks due to network speed constraints for very large swarms.

1.3 Blockchain

TODO: Write this out. The point of this section is to give the disadvantages of blockchain to support my idea to remove it

Blockchain is *TODO: Short intro of what blockchain is. This does not need to be long as this is not important* [10]

Blockchain has problems with scalability *TODO: Why is blockchain hard to scale* [11].

Blockchain is inefficient and slow *TODO: The following paper shows that blockchain can be very slow* [11].

Chapter 2

Problem and Goals

2.1 Problem

In machine learning, there exists some issues that arise from trying to use conventional techniques in the real world.

2.1.1 Privacy

It is common for data to be spread across multiple locations, referred to as data islands. Traditionally, all of this data would be consolidated into a single centralized server to facilitate the process of machine learning. However, it may not be possible to do so, given the potential conflict with privacy legislation. This leaves two options to the data scientist who is looking to train a model: either a single model per data island, likely with inferior performance, or use an algorithm that would allow the different data islands to collaborate and collectively train a model.

Consider this scenario where many different hospitals wish to train a model to detect an illness in a patient. However, due to the obligation to maintain patient confidentiality, the medical data of patients cannot be shared with any of the other hospitals. This means that, despite likely having superior performance, a model trained on all data across all hospitals is not feasible, as a consequence patients would not have the highest quality medical care possible.

2.1.2 Performance

In general, it has been observed that larger machine learning models coupled with more data typically lead to better performance. However, the use of conventional approaches for training these large models necessitates the requirement of a powerful computer. Most entities, however, do not have access to such a computer. Nevertheless, they may have access to multiple, lower-power computers. For instance, during non-working hours, a company may have hundreds of computers in its offices that are not being used, thus providing a pool of unused processing power.

2.2 Goals

2.2.1 Design an Novel Algorithm for Swarm Learning

In this project, the primary aim is to design a novel SL algorithm. This algorithm should be based on FedAvg and *Swarm Learning*. However, the algorithm should be fully decentralised, and it should operate without a

blockchain. The decision to not use a blockchain was made due to the disadvantages stated in Section 1.3.

2.2.2 Implement the Algorithm

The primary purpose of the algorithm is to test its viability when compared to other techniques, not to be as efficient as possible for real world usage. The other important feature of this implementation is that it is easy to replicate by a data scientist who wished to use SL in their next product. For this reason, when implementing the algorithm, an easy-to-understand programming language should be used, and the focus should be on readability rather than absolute performance.

2.2.3 Test Performance in Situations Where FedAvg Performs Well

FedAvg is designed to work in situations where each node has a reliable direct connection to the server. It should be shown that the new algorithm can perform well in similar situations, where each node has many stable connections to other nodes in the network.

2.2.4 Test Performance in Situations Where FedAvg Performs Badly

FedAvg may not be effective when the server has unreliable connections to the nodes, or when nodes are halted. Additionally, if the server stops working, FedAvg is unable to proceed. The goal is to demonstrate that the new algorithm can achieve successful results in scenarios where nodes frequently drop out of the network, thus making it possible to use the algorithm in cases where FedAvg may not be viable.

Chapter 3

Analysis and Specification of the Solution to the Problem

TODO: I don't know what is supposed to go here

Chapter 4

Detailed Design

DRAFT

4.1 Node

In SL, a node is an agent responsible for facilitating the improvement of the global model. The global model is an abstract concept representing the consensus of all nodes in the network. In the proposed version of SL, referred to as SwarmAvg, each node maintains its own model, known as the local model, which is an approximation of the global model. However, in SwarmAvg, the consensus algorithm used is repeated averaging, not blockchain. This means that at the start of each training step, every node may start with a slightly different model. As training progresses and performance begins to plateau, each node's local model should not only converge towards a minima, but also towards each other. Over time, each node's local model better approximates the global model, and the global model becomes a better solution to the problem.

Each node in the network possesses a confidential dataset that is not disclosed to any other nodes. In order to train the global model, nodes fit their own local model of their local dataset. In order to maintain consistency between local and global models, a combination procedure is conducted following each round of local training, which involves the integration of neighbouring nodes' models into the local model.

The steps in each training loop for SwarmAvg are as follows:

1. Fit local model to local dataset.
2. Send local model to all neighbours.
3. Combine neighbouring models into local model, using one of the methods presented in Section 4.4.

In addition, each node retains a local cache of the most recent models of its neighbouring nodes. This cache is updated each time a neighbouring node transmits its model to the node in question, instead of being updated on-demand during the combination step. The reasoning behind this decision is elaborated upon in greater detail in Section 5.3.1.

4.2 Disposal of Blockchain

TODO: Write this out

- Neural nets are heuristic - they don't need to be exact
- Its just overhead

- SL may benefit from large networks, which is bad for bc
- SL needs fast BC updates, which is hard

4.3 Training Counter

A vital aspect of SwarmAvg , specifically the combination step, involves evaluating the performance of a local model. The conventional approach would involve testing each model using an independent test set. However, due to the inability to exchange test sets among nodes, this approach is not feasible as it would result in non-comparable scores for each model. In order to circumvent this problem, this paper presents a heuristic metric referred to as the "training counter," which serves as an approximation of the level of training for a network by estimating the number of training steps performed on a given model.

The training counter can be changed in one of two manners. Firstly, the counter is incremented by 1 when the local model is trained on the local dataset, indicating that an additional step of training has been performed. Following the combination step, the training counter is also updated to reflect the combination method that was utilized. For instance, if the neighbouring models were averaged, the training counter would be updated to represent the average of all neighbouring nodes' training counters.

4.4 Model Combination Methods

The combination step is a crucial component of SwarmAvg . During this step, a node merges its local model with those of its neighbours, producing an updated estimate of the global model. This paper presents multiple methods for performing the combination step.

In the below equations, $\mu(x)$ denotes the function $mean(x)$. All models are assumed to be 1-dimensional arrays, meaning that mathematical operations such as add (+) and multiply (*) can be performed in an element-wise fashion. To achieve this constraint in the context of neural networks, a simple flattening operation is used on the weight matrices.

4.4.1 Averaging

The most rudimentary approach to combination is to compute the average model between the local model and the models of all neighbouring nodes.

$$localModel \leftarrow \mu(localModel \cup neighborModels)$$

This technique is utilized in FedAvg, which is the most simplistic form of FL. The benefit of this method is that it necessitates no hyper-parameters, which means that the data scientist needs to do less tuning. However, this attribute can also be viewed as a drawback, as it affords less flexibility in terms of customization for particular tasks.

4.4.2 Averaging With Synchronisation Rate

A more complex approach to combination is to compute the average model of all neighbours, then compute the weighted average between that model and the local model.

$$localModel \leftarrow (1 - \alpha) * localModel + \alpha * \mu(neighborModels)$$

The synchronisation rate, denoted as α , indicates the degree to which each node adjusts its local model to align with the global model. If α is set too low, each node's model in the network will diverge, resulting in each node becoming trapped at a local minima. On the other hand, if α is too high, the progress achieved by a given node will be discarded at each averaging step, which can result in slower learning.

4.4.3 Filtering By Training Counter

A potential modification to the previously mentioned combination algorithms involves filtering based on the training counter. Specifically, a node may only include its neighbour models if they meet the following statement:

$$neighborTrainingCounter + \beta \geq localTrainingCounter$$

The training offset β is the amount the training counter of a neighbour can be behind the local training counter before it is ignored.

A compelling reason to allow training counter filtering is the increased fault tolerance. Consider the situation where node A has received many model updates from many nodes, one of which being node B . However, node B goes offline and no longer is sending model updates. Without training counter filtering, node A will continue to combine the outdated node B model with it's own for as long as it is training. However, if training counter filtering is enabled, after a number of training steps the outdated model B updates will be ignored.

An issue with training counter filtering pertains to the presence of run-away nodes. These nodes possess a substantially higher training counter compared to all other nodes in the network, meaning that when filtering is

applied, they are left with no neighbours to utilise in the combination step. Consequently, a runaway node may start to overfit on its own training data, as this is the only data it is exposed to, thereby leading to decreased local performance, as well as potential performance reductions in the rest of the network. To address this problem in SwarmAvg, each node must wait until it has obtained at least γ viable neighbours prior to performing the combination step. Although this measure prevents individual nodes from becoming runaway nodes, groups of size γ still have the potential to become runaway as a unit. Nevertheless, if γ is roughly equivalent to the number of neighbours and all neighbours train at a comparable rate, the issue is minimised.

4.5 Sparse Network Behaviour

Given the sparsely connected nature of distributed scenarios, it is often the case that nodes only have direct connections to a small subset of their neighbours. This is a situation where FL struggles, however the author hypothesises that SwarmAvg should be able to deal with this situation.

4.5.1 Passive Convergence

An approach to deal with a sparsely connected network is to use the swarm learning algorithm without any modifications. This approach is effective due to the use of averaging as a combination method. When a node tries to update the global model, its changes will propagate through the network slowly, over many training iterations, even to nodes that are not directly connected. This approach has the advantage of requiring no extra data transmission, resulting in significantly less data traffic compared to other methods.

However, this method also has certain theoretical drawbacks. Consider a scenario where the network is comprised of several sparsely connected groups of nodes, where each node in a group is densely connected to other nodes within that group. In this case, it is possible that each group may learn a distinct solution to the problem. This is inefficient because instead of functioning as a cohesive network, there are multiple smaller networks acting somewhat independently of each other, potentially leading to a decrease in overall performance.

4.5.2 Relay

A solution to this could be to relay any received model updates, which means that as long as each node has at least one path to reach all other nodes, the

network will behave as a dense network. This approach offers theoretical immunity to changes in network topology, but in practice, the network's performance may still decrease compared to a truly dense network due to slower communication times between non-connected nodes.

The main disadvantage of this approach is the drastic increase in network traffic, which in turn will lead to longer model transfer times. If the swarm learning algorithm is applied to a low power network, such as an IoT network, the increase in network traffic may not be feasible at all. This relay approach can also be applied to federated learning with the same advantages and disadvantages. For these reasons, Relay will not be discussed further.

4.6 Complete Algorithm

The complete pseudocode algorithm for SwarmAvg can be found at Appendix A. It has multiple parts which are described below.

The update loop, which can be found at Appendix A.1, is the section of the algorithm which runs continuously during the time which a node is running. It takes care of training and synchronising the local model. The provided code represents what happens in a single update step, meaning that it should be run in a loop that terminates once a stop condition, such as target accuracy, has been reached.

The model received event, which can be found at Appendix A.2, is run on the local node every time a remote node sends the local node a model update. This event takes care of updating the local model cache, to ensure the local node has the most up-to-date information. The model update should contain the model of the remote node and also the remote nodes training counter.

Chapter 5

Implementation

DRAFT

5.1 Dataset and Machine Learning Model

5.1.1 Dataset

Initially, the dataset utilized for experimentation was the MNIST dataset, which encompasses 60000 greyscale 28x28 labelled images of digits from 0 to 9. This dataset was selected for its simplicity, requiring no pre-processing or data cleaning prior to training, and due to its availability as a built-in component of the chosen machine learning framework, Keras.

However, upon implementation it was determined that this dataset was too simple for the application of machine learning, as a single node could reach near peak accuracy after a single epoch, rendering it ill-suited for swarm learning, an algorithm designed to function across multiple training epochs.

To address this issue, the MNIST dataset was replaced with MNIST-fashion, a drop-in replacement dataset containing 10 classes of different items of clothing. MNIST-fashion is known to be more challenging [12]. To further increase the complexity of the problem, in several experiments each agent was only provided with a small subset of the entire dataset, resulting in less training per epoch, and therefore meaning that an agent would require more epochs to achieve the same performance.

5.1.2 Model

The Keras machine learning framework in Python was utilized to implement the model due to its reputation for being both simple and straightforward. All experiments made use of the same model, which is outlined in Appendix B and is a small convolutional neural network. The model was tested on the MNIST fashion dataset and was able to attain an accuracy score of above 90 percent when trained on the entire dataset using local learning; this result is on par with the accuracy reported in the original paper [12], making the model suitable for use.

5.2 Federated Learning

FedAvg was chosen for comparison of performance against SwarmAvg . In order to ensure fairness of the comparison, it was necessary to implement FedAvg from scratch using the same language framework as SwarmAvg.

5.2.1 Algorithm

The implemented algorithm worked in the same manner as the algorithm described in the FedAvg paper. However, one modification was made: at the start of each timestep instead of choosing N random nodes to perform training, all nodes were chosen. This means that every available node will perform training at every timestep, which should result in the best possible performance for federated learning, especially given that only 10 nodes could be run at once.

Initially, a REST API was utilised to transfer the model between the server and the client. However, this was deemed unnecessary and was supplanted for two reasons. Firstly, the decision was made to measure performance against epochs trained instead of performance against time, which is discussed in further detail in the results section. Secondly, the REST API approach was much slower than the method chosen to replace it, yet it resulted in the same performance measurements when gauged in terms of epochs trained. As a substitute for the REST API, a system of functions was implemented. When a node sent a model to another node, it simply called a function on that node. This was abstracted away from the main algorithm code, thus meaning that a drop-in REST replacement could be added at a later date. This significantly accelerated training, and enabled a greater number of training runs to be conducted, resulting in more data being collected.

5.2.2 Evaluation

This implementation of federated learning is simple yet effective. Though certain simplifications have been made, they should not interfere with the performance of the model in this scenario. It should only be used for testing performance against epochs trained, not against time taken, as the model transfer layer has been simplified.

5.3 Prototype

It was decided to make an initial, less streamlined, prototype as a proof-of-concept before spending a lot of time creating the final algorithm. This prototype was created to be very modifiable so that changes could easily be made and tested.

5.3.1 Algorithm

This algorithm was a simplified version of that described in the design section. The primary difference was that when a node performed its synchronisation step, it would request the models from each neighbour instead of utilising its cached versions of their models. This had the consequence of slowing down training, as the synchronisation step could not be completed until all nodes had responded. Furthermore, this algorithm did not incorporate the training counter, leading to the absence of training counter filtering, β and γ .

5.3.2 Evaluation

This step was beneficial for the progression of the project, as it enabled the author to form the ideas detailed in the design process, which were then implemented in the subsequent step. However, due to the abundance of superfluous code, the algorithm was inefficient and performed poorly. For this reason, it was decided not to record the results of this method.

5.4 Final

In this implementation, the findings of the prototype were taken and built upon. The code was streamlined for the purpose of testing performance. However, the code was still designed to be reusable and easy to read.

5.4.1 Algorithm

The algorithm is as described in the design section. However, there was one discrepancy which needed to be tested: whether to send model updates to neighbours before or after performing the combination step. Both were tested but pushing model updates before combination seemed to be the more effective method when comparing the accuracy. The author hypothesises that this is due to more of the local nodes training progress being preserved, meaning that a local nodes training has more of an affect on its neighbours.

The back end for distributing models was implemented as an interface which abstracts the details of the distribution away from the main algorithm code. For this implementation, the same distribution strategy as the previously implemented FedAvg was used: calling local functions that simulated a web connection.

5.4.2 Evaluation

Overall, this implementation of the proposed swarm learning method is satisfactory. It is efficient, reusable and simple to understand and use. Nevertheless, it is not yet suitable for real world applications, as it lacks any security features and there is limited error handling, being designed for use in a controlled testing environment. However, it would not be challenging to incorporate a backend for this code, allowing for communication over the internet.

Chapter 6

Testing Strategy and Results

6.1 Methods

6.1.1 Metrics

The metric chosen for the following experiments was accuracy, due to its comprehensible nature. Despite the fact that an unbalanced dataset presents one of the significant drawbacks of using accuracy, this concern is irrelevant in the case of MNIST since it is a balanced dataset. The accuracy of each node is computed after each training step using the test subset of MNIST, and none of the nodes are ever provided access to the test set for training.

6.1.2 Data Collection

The training process for each experiment was conducted five times, and the resulting accuracies of every node were recorded. To mitigate the impact of training noise on the performance graphs, the accuracy value for each time step was calculated as the median accuracy across all nodes and runs at that time step.

6.1.3 Node Counts

The experiments were conducted using 10 nodes, with the exception of the server in cases where FL was employed. The decision of how many nodes to simulate was based on the highest node count attainable without causing inconsistencies and crashes due to resource depletion of the training machine.

6.1.4 Algorithm Configurations

In each of the following experiments, the algorithm was configured using a specific set of parameters $\alpha\beta\gamma$. These parameters were obtained heuristically by making an initial guess, testing, and then fine-tuning them until a satisfactory outcome was reached. However, it is important to note that an exhaustive investigation into the optimal parameter configuration for a particular type of problem is not within this paper's scope, meaning that it is plausible that swarm learning could yield better results with more precisely tuned parameters.

6.1.5 Data Volume Per Node

The experiments evaluate the algorithm's performance using three levels of data volume per node. These levels are considerably smaller than the full

MNIST dataset not only to increase problem difficulty, but also as the algorithm is intended for scenarios where each nodes access to data is restricted. To create a subsection of data for each node, a random sampling with replacement method was used to select the desired number of datapoints. During the initial training phase, each node performs a single sampling of its dataset, after which that nodes data subset remains constant.

6.1.6 Epochs

Due to the limited size of the dataset, a single node executes more than one epoch of training in each training loop. The number of epochs carried out by a node per training step will be referred to as Epochs per Step (EPS). Empirical testing has indicated that both SL and FL exhibit improved performance with higher EPS, at times surpassing the gains from increasing the number of training steps. Moreover, the utilization of higher EPS was favoured due to its reduced training time, compared to increasing the number of training steps. The three levels of data volume with their respective EPS are shown in table 6.1

Dataset Size	EPS	Reason
6000	2	As there are 10 nodes, it was decided that each node should be tested with 1/10th of the dataset
1000	5	A much lower volume of data was tested to reflect the anticipated use case of SL - training models where each node has very small amounts of data
100	15	The extreme case was tested to see how well the algorithms perform in undesirable conditions

Table 6.1: The different levels of dataset size and EPS that were tested

6.2 Dense Network Performance

A crucial experiment for evaluating the performance of the SL algorithm involves assessing its performance under optimal circumstances, specifically within a network of nodes wherein each node is directly connected to every other node. In FL, the analogous topology involves direct connections between each node and the server. This comparison is significant as it facilitates a direct evaluation of the SL and FL algorithms under their respective ideal conditions.

The selection of the parameters for the SL algorithm was based on the authors prior experience in testing the algorithm. These parameters are presented in Table 6.2.

P	Value	Reason
α	0.75	Low enough to allow nodes to maintain a small variation but not so low that the nodes diverge indefinitely
β	0.5	Allows a small amount of nodes looking back, but high value is not needed as every node will be running at approximately the same speed.
γ	8	All nodes will always be connected to 9 other nodes, so higher is better. There is room for 1 node to be skipped to prevent deadlock.

Table 6.2: The chosen parameters for the SL algorithm

6.2.1 Results

Below are the results attained from running the training script. In each graph, SwarmAvg is in red and FedAvg is in black. The shaded area around each line is the upper and lower quartile.

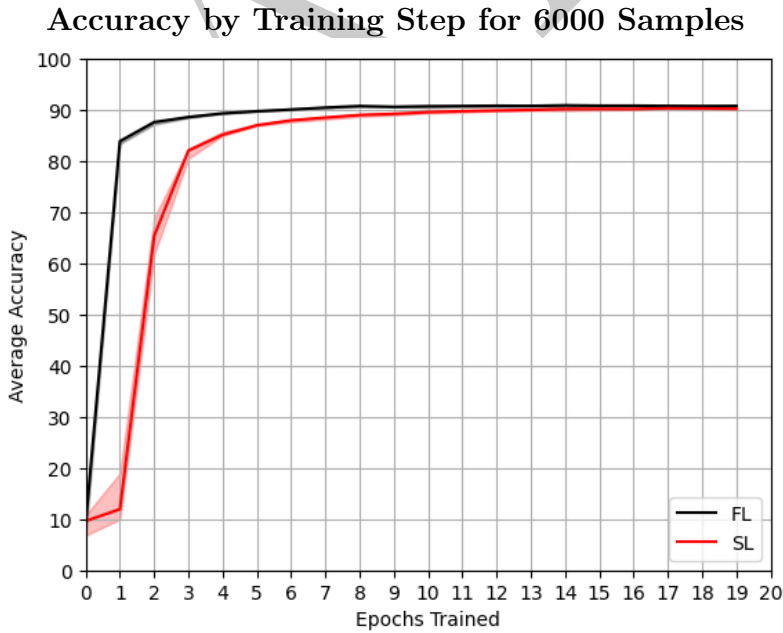


Figure 6.1: Comparing Accuracy of FL and SL with 6000 Data Samples per Node

In Figure 6.2, it is clear that SwarmAvg is the slower of the two algorithms. However, in this situation both algorithms reach approximately the same final accuracy. The main difference is that the performance of SwarmAvg lags behind FedAvg by 1-2 epochs.

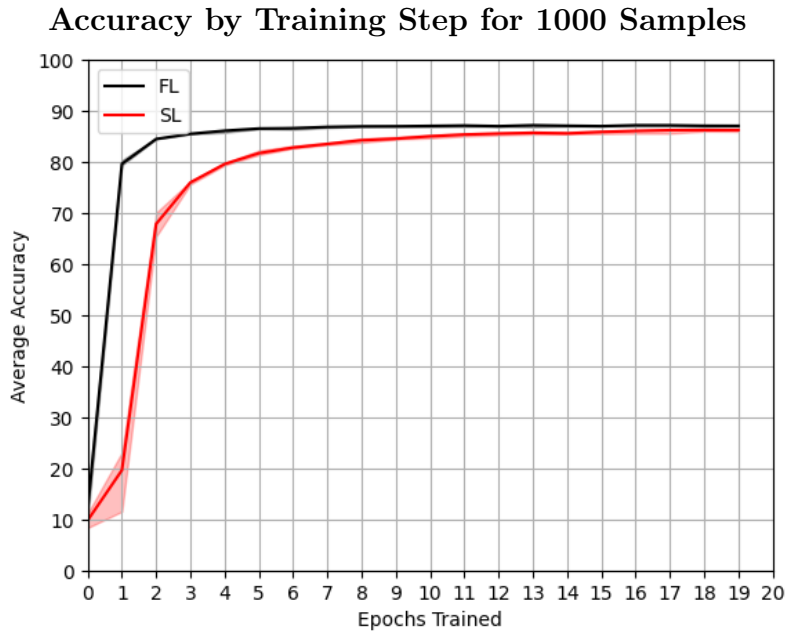


Figure 6.2: Comparing Accuracy of FL and SL with 1000 Data Samples per Node

Unlike Figure 6.1, Figure 6.2 clearly shows that SwarmAvg finished training with a lower accuracy. SwarmAvg also still lags by approximately the same amount. Due to the reduction in data, both techniques have a lower accuracy than they previously achieved.

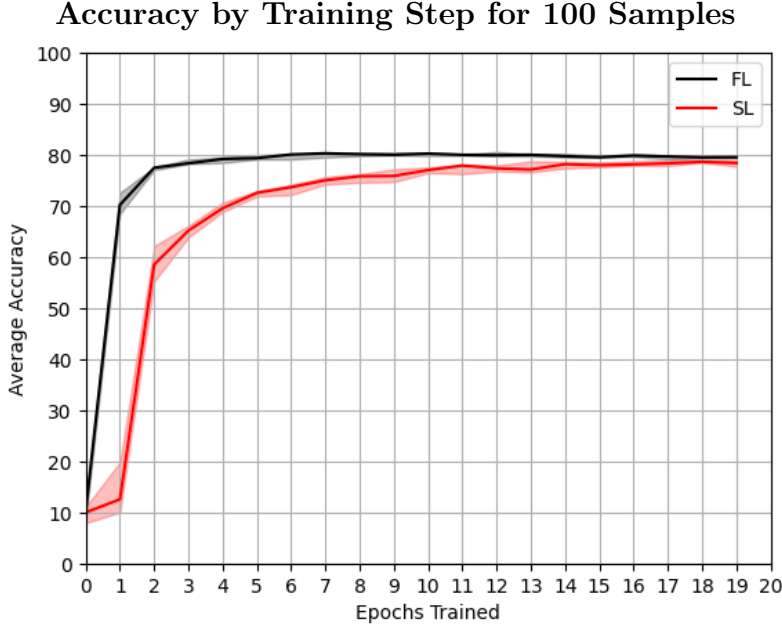


Figure 6.3: Comparing Accuracy of FL and SL with 100 Data Samples per Node

The trend of SwarmAvg ending with a lower accuracy than FedAvg continues in Figure 6.3. Both algorithms also have a much lower accuracy than what they achieved previously, due to the drastic decrease in available data.

6.2.2 Analysis

The most evident change when reducing the volume of data is the large reduction in accuracy of the two algorithms, which was an expected observation. However, It can also be observed that SwarmAvg seems to be affected by this reduction in data slightly more than FedAvg is, due to its steady increase in difference in peak accuracy. Despite this, the difference in peak accuracy is very small, often within 2 percent, which shows that SwarmAvg does appear to be working correctly.

One of the major problems with SwarmAvg is its consistently slower convergence on an optima than FedAvg. Especially from the onset, SwarmAvg has a much lower accuracy than FedAvg, and takes longer to reach its peak accuracy. This may be attributed to the asynchronous manner of the nodes in SwarmAvg, meaning that some nodes that perform training before all others have a lower accuracy than what would be expected.

It should also be noted that SwarmAvg has a higher inter-quartile range than FedAvg in all of these tests, meaning that FedAvg is more consistent

with respect to accuracy. However, this difference is marginal.

TODO: Overfitting?

6.3 Dense Network Performance with Node Dropout

TODO: Do this experiment - it really shouldn't take long

the same but this time dropout nodes. Try with and without filtering to show that it increases fault tolerance

To do this, test where a certain number of nodes drop out at step 1, 2, 3, etc

6.4 Sparse Network Performance

In reality, it is uncommon for each node to be linked with every other node. To simulate a more realistic scenario, a technique was employed to generate a network of nodes with a specific density. It is important to note that, moving forward, density refers to an artificial metric and is not associated with the physical definition of density. When density is set to 0, the network is minimally linked, meaning that each node has at least one indirect path to every other node, but the minimal number of connections required to accomplish this exist. When density is set to 1, all nodes are connected to one another in a dense fashion. Since this measure may not provide a straightforward indicator of network density, two additional metrics will be provided: Mean Minimum Hops (MMH) and Mean Connections per Node (MCPN). MMH denotes the mean minimum number of transitions required to get from one node to another in the network. MCPN denotes the average number of connections a given node possess.

To ensure good test coverage for a multitude of possible deployment conditions, multiple densities were tested for each count of data. Table 6.3 lists the densities that were tested along with their corresponding MMH and MCPN.

Density	MMH	MCPN
1	1.0	9.0
0.75	1.2	7.2
0.5	1.4	5.4
0.25	1.7	3.6
0	3.0	1.8

Table 6.3: The statistics for different density levels

To help the reader to visualise the different density levels, Figure 6.4 shows a visual representation of an example network that was generated for each density.

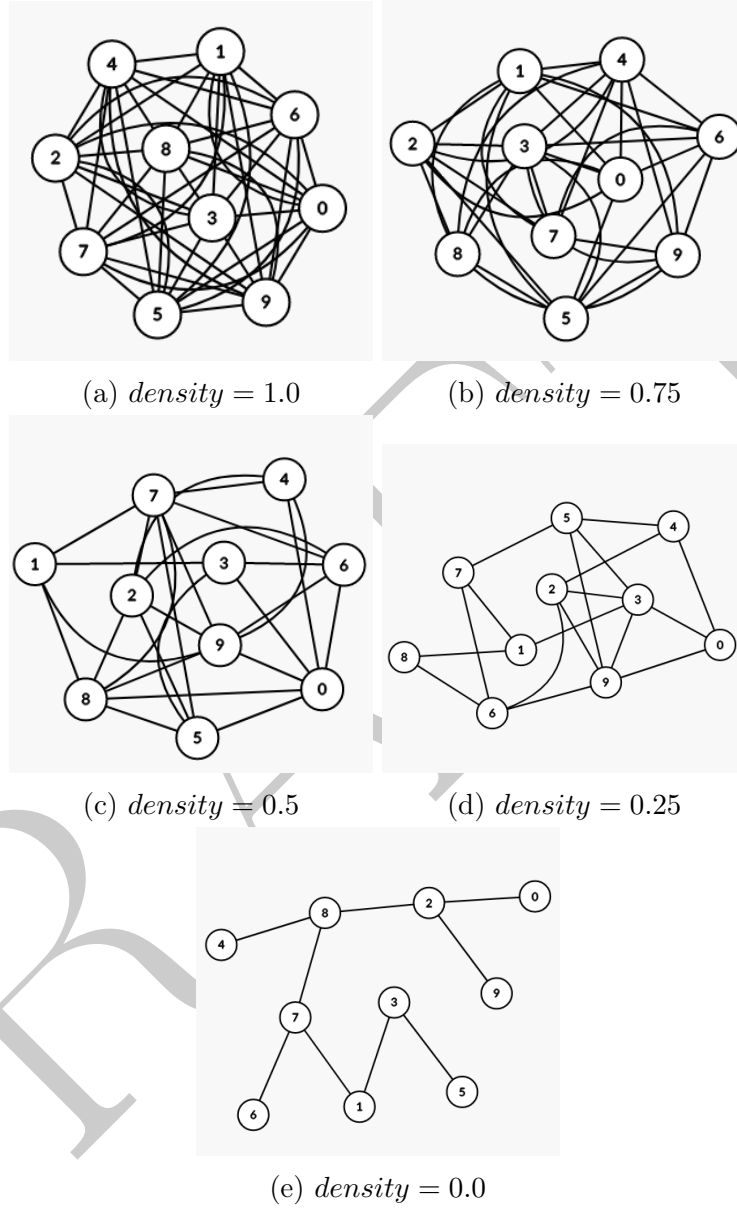


Figure 6.4: Example networks of nodes generated for each density level, visualised using the tool at [13]. Each time a simulation is started, a new random network is generated for that simulation.

In this section, FedAvg is not tested. When a every node is not connected

to the server node, FedAvg has two options. It can either ignore all nodes that are not directly connected to it, or it can try to relay the model updates. For the reasons mentioned in Section 4.5.2, relaying is not always a good solution, so the decision was made to only test SwarmAvg in this test.

6.4.1 Results

Below are the results attained from running the testing script. The more green a line is on the graph, the higher the density, and the more red, the lower the density. For each of the tested densities, γ was set to $\text{round}_{\text{down}}(\text{MCPN}) - 1$, meaning that each node would usually wait for all but one of its neighbours until it could progress. If γ was not reduced for the less dense networks, many of the nodes would end up waiting for a number of neighbours that was not possible to achieve.

Accuracy by Training Step for 6000 Samples for Different Densities

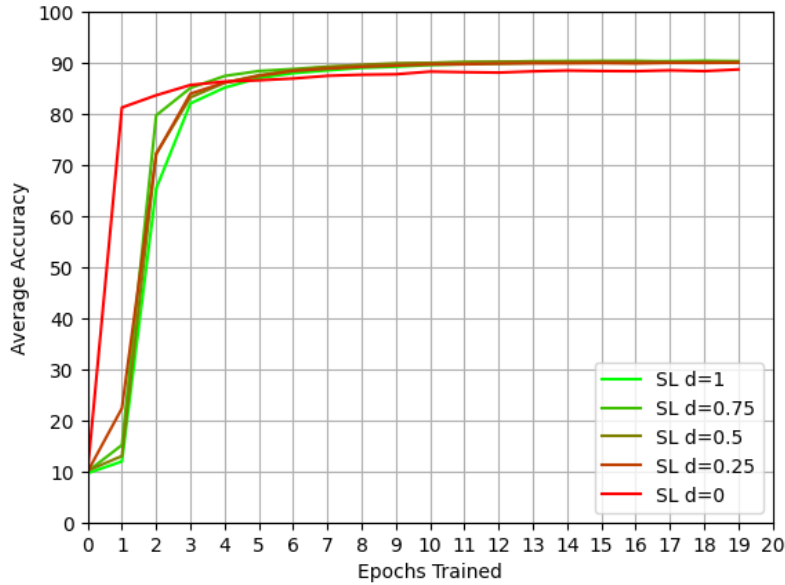


Figure 6.5: Comparing Accuracy of SL with 6000 Data Samples per Node and varying network density

Interestingly, with 6000 samples, all the tested densities except for 0 seem to perform with very similar performances. They all reach the same peak accuracy and converge in a similar amount of time. At density 0, convergence is actually faster but the final accuracy is lower overall.

Accuracy by Training Step for 1000 Samples for Different Densities

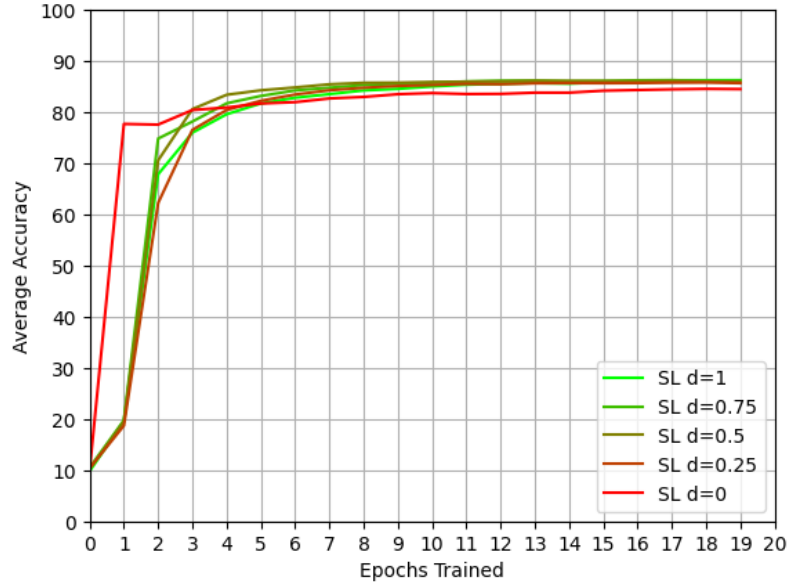


Figure 6.6: Comparing Accuracy of SL with 1000 Data Samples per Node and varying network density

With 1000 samples, the tests yielded relatively similar results to 6000 samples, except for the fact that the overall accuracy for all tests was lower.

Accuracy by Training Step for 100 Samples for Different Densities

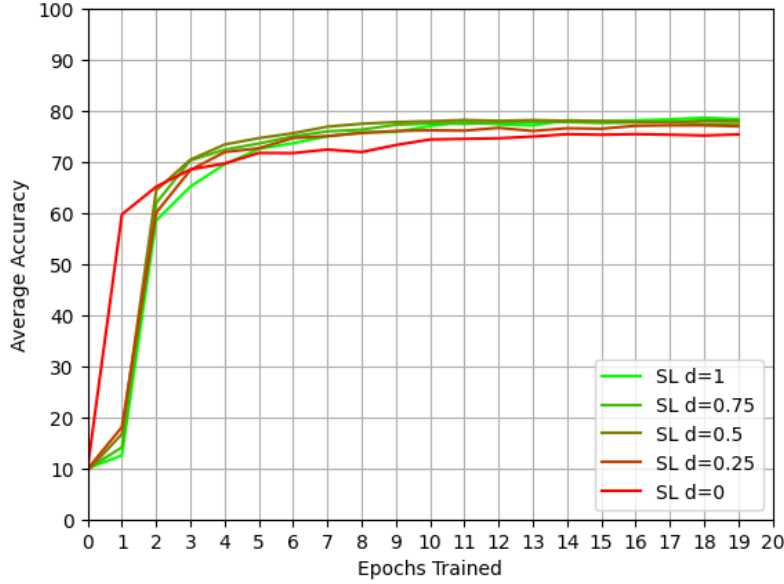


Figure 6.7: Comparing Accuracy of SL with 100 Data Samples per Node and varying network density

At just 100 samples, it is clear that the density of the network is affecting the final accuracy. The final accuracies range from between 75 and 80 percent, where the more dense the network, the higher the final accuracy. The training curves from Figure 6.7 also seem to be more noisy than the tests with higher volumes of data.

6.4.2 Analysis

In general, changing the density of the network of nodes has a small affect on the training of the nodes in that network. A lower density tends to converge slightly faster but tend towards a lower final accuracy. This effect is more pronounced as the volume of data provided to each node reduces, which is demonstrated by the higher variation in training shown in Figure 6.7 when compared to Figure 6.5 and 6.6.

The networks with density of 0 converge towards their optima faster than the higher density networks. In figure 6.7 it can also be seen that the lower densities converge slightly faster than the higher densities but are over taken by the end of training.

6.5 SwarmAvg on a FedAvg topology

TODO: Im not sure if this is a bad thing to include, but it might be interesting to say that if you want to switch to swarm, but already have a centralised network, this is the performance to expect

Chapter 7

Critical Evaluation

7.1 Comparison of SL to FL

TODO: Make this sentence into a paragraph with some more detail?

This report has show that SL has the capability to be a competitor to FL in a distributed learning situation. SL can also be used in some situations where FL cannot. However, there are some clear disadvantages to using SL as opposed to FL.

7.1.1 Advantages of using SL over FL

TODO: Write this out

- Can prevent a single slow node from bottlenecking the whole process
- No central server so more fault tolerant
- Can still function if all nodes do not have direct connections to all other nodes (in the case of FL every node must have access to the server)

7.1.2 Disadvantages of using SL over FL

TODO: Write this out

- Less data transmission ($O(n)$ rather than worst case $O(n^2)$). This effect is less prominent if SL is sparse
- FL is a Less complex algorithm -> fewer parameters to tune -> easier to fit to a problem
- Slower to converge

7.1.3 Limitations of the Presented Implementation

can diverge, unlike using blockchain

needs an extra layer of security, but this is built into blockchain

7.1.4 Further Extensions of Presented Implementation

could be used for any situation where a group of agents need to come to a consensus of a moving targeted thing, where each agent updates the target slightly

Chapter 8

Project Management

8.1 Time Management

To assist with planning and organisation of this project, Gantt charts were used. These helped to visualise which tasks needed doing, and also if the project progress was ahead or behind schedule. The Gantt for the interim report can be found in appendix C.1 and the final report Gantt can be found at C.2.

On the Gantt charts, it was decided to not count weekends as working days. This was chosen as it represented the fact that the author had other work to do during the course of the project.

The presented charts represent the final iteration of planning. As described in the following section, changes occurred throughout the project and the Gantt charts were adapted accordingly.

8.2 Changing Plans

At the time of writing the interim report, the project plan was very different to what is presented in this report. Originally, the plan was to implement SL and then to simulate real world problems such as uneven data distribution and sparse networks. However, during the process of implementing the SL algorithm it was apparent that the SL algorithm itself would take a very large amount of time to develop and test. For this reason, the decision was made to refocus the project onto the development of the SL algorithm, and if enough time was available after this, testing a single real world problem may be possible.

Plans were also changed on a smaller scale somewhat regularly. This is due to the experimental nature of this project, and the fact that it was impossible to predict some problems during the planning phase. For example, the parameters β and γ were not planned for at the start of the project, but whilst experimenting problems arose which seemed like they could be fixed by those parameters. This meant that another implementation had to be created which included those parameters.

The small changes to the plan were not an issue due to the inclusion of many buffer zones in the planned time allocation, such as *Bugfixing Time*. This meant that an overrun task would not affect too many tasks moving forwards, as it could be caught up on in a buffer zone.

8.3 Risk Assessment

The risk assessment was created early on in the project to mitigate any of the major risks. However, none of the described problems arose to an extent which their mitigation plan needed to be followed.

8.3.1 Personal Issues

Description

Personal issues which cause the author to be unable to do work, such as illness.

Risk Calculations

Severity (1-5): 3

Likelihood (1-5): 3

Overall Risk (1-25): 9

Mitigation

The codebase will be designed such that individual sections and modules have minimal dependencies on other sections. This means that, even if the author is unable to work for a period of time, some less critical sections can be omitted without significantly impacting the rest of the project.

8.3.2 Hardware Failure

Description

Failure on the authors local computer of any kind, such as a graphics card or storage breakage.

Risk Calculations

Severity (1-5): 4

Likelihood (1-5): 2

Overall Risk (1-25): 8

Mitigation

The project will be regularly backed up to GitHub. If a core component of the work computer breaks, the author has access to a personal laptop and

the Zepler Labs. The deep learning environment along with dependencies is backed up to the authors Google Drive in the form of a docker image, so that switching to a new computer would be a smooth process.

8.3.3 Algorithm Does Not Work Work

Description

The SL algorithm does not function as well as expected.

Risk Calculations

Severity (1-5): 5

Likelihood (1-5): 1

Overall Risk (1-25): 5

Mitigation

It may be possible to shift the project away from SL and onto distributed FL with leader election. FL is more commonly used and therefore has more literature, meaning that it is more likely to be an achievable goal to implement it.

Chapter 9

Conclusion and Future Work

9.1 Conclusion

TODO: Conclude what is SL

TODO: Conclude the findings + advantages

TODO: Conclude the disadvantages

9.2 Future Work

TODO: Write this out

- Exhaustive search on optimal parameters for different situations
- Simulate more nodes
- Try algo out distributed amongst many machines over the internet

Appendix A

Sword Learning Algorithm

A.1 Training Step

Algorithm 1 Training Step - Called Repeatedly in a Loop

```
1: TRAIN(localModel, localData)
2: for all  $n \in neighbors$  do
3:   SENDTO((localModel, localTrainingCounter), n)
4: end for
5: for  $x \in range(maxSyncWaits)$  do
6:    $neighborModels \leftarrow \emptyset$ 
7:   for all  $n \in neighbors$  do
8:      $model, trainingCounter \leftarrow CACHELOOKUP(n)$ 
9:     if  $trainingCounter + \beta \geq localTrainingCounter$  then
10:      APPEND( $neighborModels$ ,  $model$ )
11:    end if
12:  end for
13:  if  $length_{neighborModels} \geq \gamma$  then
14:    if  $synchronisationMethod = "AVG"$  then
15:       $localModel \leftarrow \mu(localModel \cup neighborModels)$ 
16:    else if  $synchronisationMethod = "ASR"$  then
17:       $localModel \leftarrow (1 - \alpha) * localModel + \alpha * \mu(neighborModels)$ 
18:    end if
19:  else
20:    continue
21:  end if
22:  SLEEP(syncWaitTime)
23: end for
```

A.2 Model Received Event

Algorithm 2 Model Received Event - Called When a Model Update is Received from a Remote Node

Input: *neighbour, nModel, nTrainingCounter*

```
1: if INCACHE(neighbour) then
2:    $\_, nTraningCounterOld \leftarrow \text{CACHELOOKUP}(n)$ 
3:   if nTrainingCounter > nTraningCounterOld then
4:     SETCACHE(neighbour, (nModel, nTrainingCounter))
5:   end if
6: else
7:   SETCACHE(neighbour, (nModel, nTrainingCounter))
8: end if
```

Appendix B

Machine Learning Model

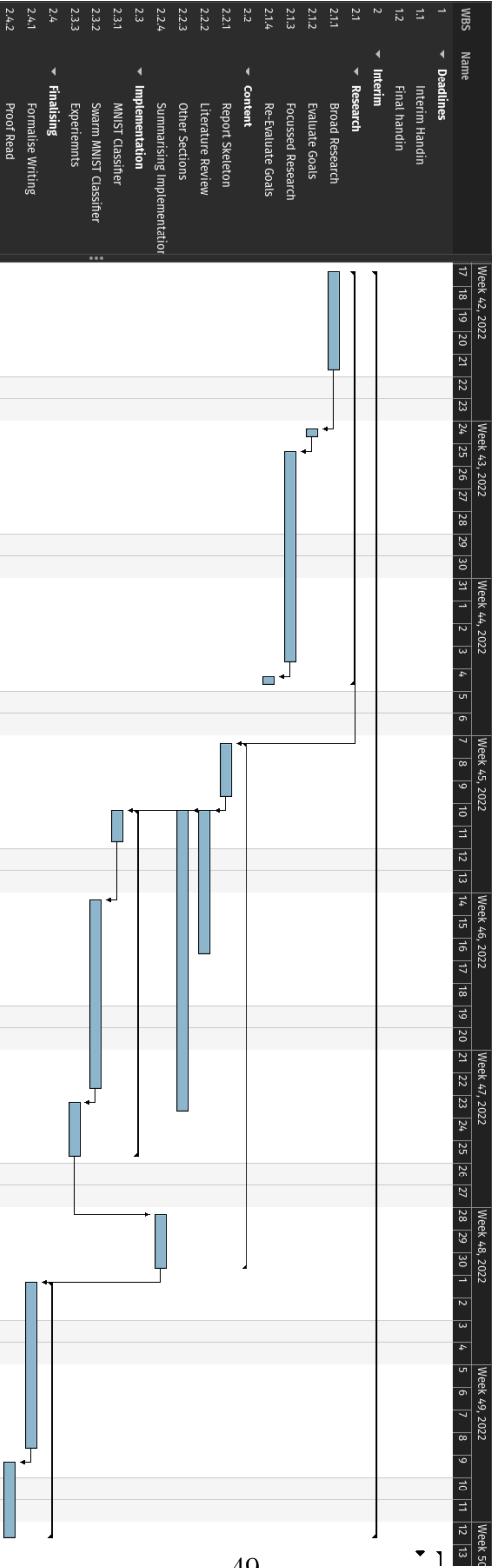
```
inp = Input((28,28))
out = Reshape((28,28,1))(inp)
out = Conv2D(16, (3,3), activation="relu")(out)
out = Conv2D(16, (3,3), activation="relu")(out)
out = Flatten()(out)
out = Dense(256, activation="relu")(out)
out = Dense(128, activation="relu")(out)
out = Dense(10, activation="sigmoid")(out)
model = Model(inputs=inp, outputs=out)
model.compile(
    optimizer="adam",
    loss=SparseCategoricalCrossentropy(),
    metrics=[SparseCategoricalAccuracy()]
)
```


Appendix C

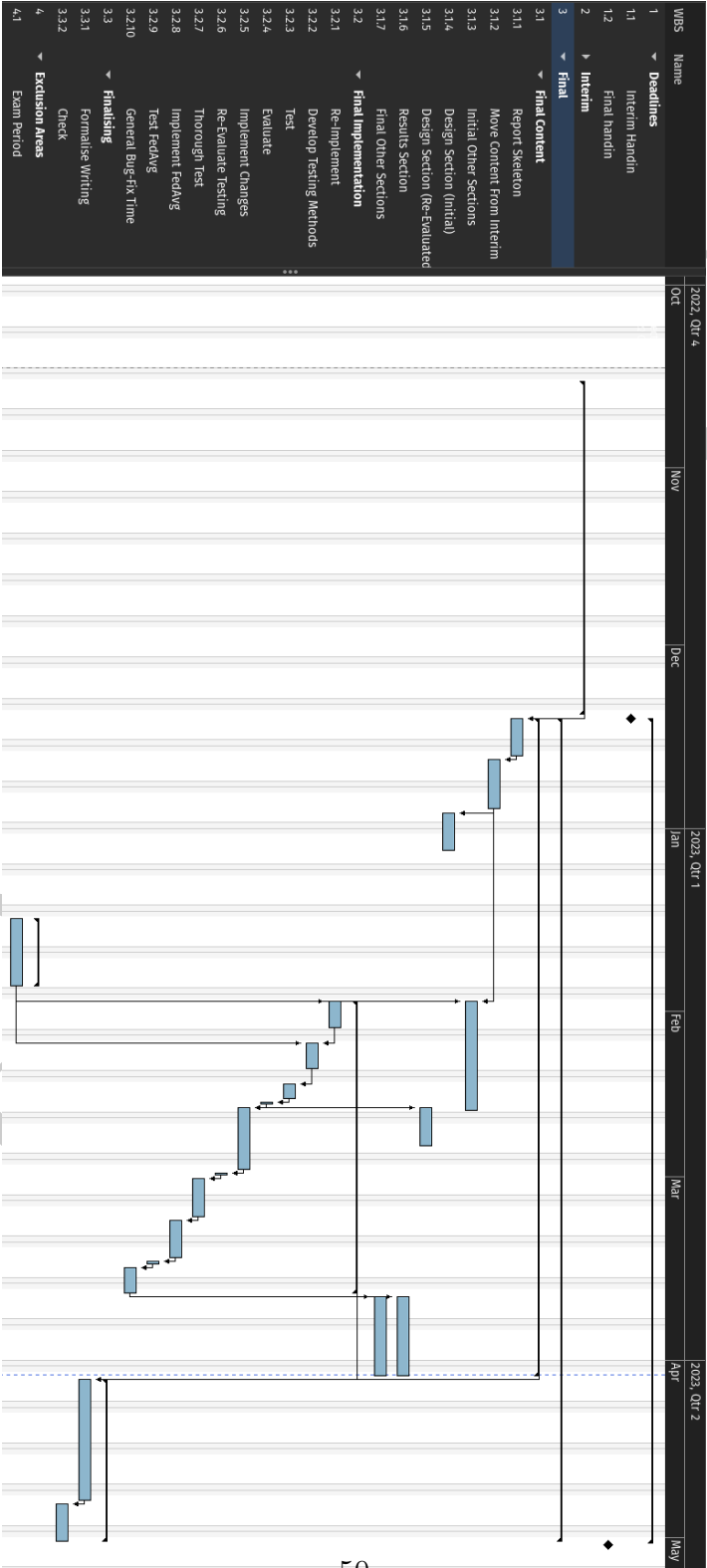
Gantt Charts

DRAFT

C.1 Gantt - Interim



C.2 Gantt - Final



Bibliography

- [1] M. Kop, “Machine learning & eu data sharing practices,” Stanford-Vienna Transatlantic Technology Law Forum, Transatlantic Antitrust ..., 2020.
- [2] G. A. Kaissis, M. R. Makowski, D. Rückert, and R. F. Braren, “Secure, privacy-preserving and federated machine learning in medical imaging,” *Nature Machine Intelligence*, vol. 2, no. 6, pp. 305–311, 2020.
- [3] C. Zhang, Y. Xie, H. Bai, B. Yu, W. Li, and Y. Gao, “A survey on federated learning,” *Knowledge-Based Systems*, vol. 216, p. 106775, 2021.
- [4] Q. Li, Z. Wen, Z. Wu, S. Hu, N. Wang, Y. Li, X. Liu, and B. He, “A survey on federated learning systems: vision, hype and reality for data privacy and protection,” *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [5] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, “Communication-efficient learning of deep networks from decentralized data,” 2016.
- [6] J.-H. Chen, M.-R. Chen, G.-Q. Zeng, and J.-S. Weng, “Bdfl: a byzantine-fault-tolerance decentralized federated learning method for autonomous vehicle,” *IEEE Transactions on Vehicular Technology*, vol. 70, no. 9, pp. 8639–8652, 2021.
- [7] S. Warnat-Herresthal, H. Schultze, K. L. Shastri, S. Manamohan, S. Mukherjee, V. Garg, R. Sarveswara, K. Händler, P. Pickkers, N. A. Aziz, S. Ktena, F. Tran, M. Bitzer, S. Ossowski, N. Casadei, C. Herr, D. Petersheim, U. Behrends, F. Kern, T. Fehlmann, P. Schommers, C. Lehmann, M. Augustin, J. Rybníček, J. Altmüller, N. Mishra, J. P. Bernardes, B. Krämer, L. Bonaguro, J. Schulte-Schrepping, E. De Domenico, C. Siever, M. Kraut, M. Desai, B. Monnet, M. Sridhara, C. M. Siegel, A. Drews, M. Nuesch-Germano, H. Theis, J. Heyckendorf, S. Schreiber, S. Kim-Hellmuth, P. Balfanz, T. Eggermann,

- P. Boor, R. Hausmann, H. Kuhn, S. Isfort, J. C. Stingl, G. Schmalzing, C. K. Kuhl, R. Röhrig, G. Marx, S. Uhlig, E. Dahl, D. Müller-Wieland, M. Dreher, N. Marx, J. Nattermann, D. Skowasch, I. Kurth, A. Keller, R. Bals, P. Nürnberg, O. Rieß, P. Rosenstiel, M. G. Netea, F. Theis, S. Mukherjee, M. Backes, A. C. Aschenbrenner, T. Ulas, A. Angelov, A. Bartholomäus, A. Becker, D. Bezdan, C. Blumert, E. Bonifacio, P. Bork, B. Boyke, H. Blum, T. Clavel, M. Colome-Tatche, M. Cornberg, I. A. De La Rosa Velázquez, A. Diefenbach, A. Dilthey, N. Fischer, K. Förstner, S. Franzenburg, J.-S. Frick, G. Gabernet, J. Gagneur, T. Ganzenmueller, M. Gauder, J. Geißert, A. Goesmann, S. Göpel, A. Grundhoff, H. Grundmann, T. Hain, F. Hanses, U. Hehr, A. Heimbach, M. Hoeper, F. Horn, D. Hübschmann, M. Hummel, T. Iftner, A. Iftner, T. Illig, S. Janssen, J. Kalinowski, R. Kallies, B. Kehr, O. T. Keppler, C. Klein, M. Knop, O. Kohlbacher, K. Köhrer, J. Korbel, P. G. Kremsner, D. Kühnert, M. Landthaler, Y. Li, K. U. Ludwig, O. Makarewicz, M. Marz, A. C. McHardy, C. Mertes, M. Münchhoff, S. Nahnsen, M. Nöthen, F. Ntoumi, J. Overmann, S. Peter, K. Pfeffer, I. Pink, A. R. Poetsch, U. Protzer, A. Pühler, N. Rajewsky, M. Ralser, K. Reiche, S. Ripke, U. N. da Rocha, A.-E. Saliba, L. E. Sander, B. Sawitzki, S. Scheithauer, P. Schiffer, J. Schmid-Burgk, W. Schneider, E.-C. Schulte, A. Sczyrba, M. L. Sharaf, Y. Singh, M. Sonnabend, O. Stegle, J. Stoye, J. Vehreschild, T. P. Velavan, J. Vogel, S. Volland, M. von Kleist, A. Walker, J. Walter, D. Wiczorek, S. Winkler, J. Ziebuhr, M. M. B. Breteler, E. J. Giamarellos-Bourboulis, M. Kox, M. Becker, S. Cheran, M. S. Woodacre, E. L. Goh, J. L. Schultze, C.-. A. S. (COVAS), and D. C.-. O. I. (DeCOI), “Swarm learning for decentralized and confidential clinical machine learning,” *Nature*, vol. 594, pp. 265–270, Jun 2021.
- [8] J. C. Varughese, R. Thenius, T. Schmickl, and F. Wotawa, “Quantification and analysis of the resilience of two swarm intelligent algorithms,” in *GCAI*, pp. 148–161, 2017.
- [9] J. Augustine, T. Kulkarni, and S. Sivasubramaniam, “Leader election in sparse dynamic networks with churn,” in *2015 IEEE International Parallel and Distributed Processing Symposium*, pp. 347–356, IEEE, 2015.
- [10] D. Yaga, P. Mell, N. Roby, and K. Scarfone, “Blockchain technology overview,” tech. rep., oct 2018.
- [11] D. Yang, C. Long, H. Xu, and S. Peng, “A review on scalability of blockchain,” in *Proceedings of the 2020 The 2nd International Con-*

ference on Blockchain Technology, ICBCT'20, (New York, NY, USA),
p. 1–6, Association for Computing Machinery, 2020.

- [12] H. Xiao, K. Rasul, and R. Vollgraf, “Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms,” 2017.
- [13] “Graph editor.”