
Consensus learning: A novel decentralised ensemble learning paradigm

Horia Magureanu
Flare Research
horia@flare.foundation

Nairi Usher
Flare Research
nairi@flare.foundation

Abstract

The widespread adoption of large-scale machine learning models in recent years highlights the need for distributed computing for efficiency and scalability. This work introduces a novel distributed machine learning paradigm – *consensus learning* – which combines classical ensemble methods with consensus protocols deployed in peer-to-peer systems. These algorithms consist of two phases: first, participants develop their models and submit predictions for any new data inputs; second, the individual predictions are used as inputs for a communication phase, which is governed by a consensus protocol. Consensus learning ensures user data privacy, while also inheriting the safety measures against Byzantine attacks from the underlying consensus mechanism. We provide a detailed theoretical analysis for a particular consensus protocol and compare the performance of the consensus learning ensemble with centralised ensemble learning algorithms. The discussion is supplemented by various numerical simulations, which describe the robustness of the algorithms against Byzantine participants.

Contents

1	Introduction	2
1.1	Main contributions	3
1.2	Related works	3
1.3	Organisation	4
2	Preliminaries	5
2.1	Ensemble learning	5
2.2	Jury problems	6
2.3	Consensus mechanisms	7
3	Consensus learning	9
3.1	Algorithm description	9
3.2	Summary of key results	10
4	Theoretical analysis	12
4.1	Modelling premises	12
4.2	Homogeneous case	13
4.3	Diversifying the base learners	18
4.4	Byzantine tolerance in consensus learning	19

5	Numerical simulations	21
5.1	Non-IID MNIST dataset	21
5.2	Beta-distributed base learners	24
6	Conclusions and outlook	26
A	Snow protocols	27
A.1	Slush with honest participants	27
A.2	Byzantine participants	30
B	Proofs of main results	31
B.1	Proof of Theorem 3	31
B.2	Proof of Theorem 1	31
B.3	Proof of Theorem 4	33
C	Simulation details	34

1 Introduction

Machine learning (ML) has traditionally existed within the context of centralised computing, whereby both data processing and computations occur on a single server. More recently, distributed settings [1–5] have garnered increased attention due the complexity of modern foundation models, such as large language and computer vision models [6], which require vast quantities of data to be processed. There, both data and computational resources can be spread across multiple devices or nodes. A prominent distributed learning paradigm is federated learning (FL), where nodes train a model collectively by sharing only local model updates in order to protect data privacy [1, 2, 7].

Yet, FL algorithms and more generally distributed algorithms are vulnerable to malicious or faulty behaviour – termed Byzantine behaviour – of the participants [8], and dealing with such participants is one of the most challenging problems in distributed ML [9–12]. The resilience against such Byzantine actors is rooted in the aggregation rule used to combine the local model updates shared in every federated training round [9]. Furthermore, even though FL algorithms can leverage encryption or differential privacy mechanisms, they remain susceptible to privacy threats [13], as the local model updates have been shown to contain enough information to reconstruct local data samples [14–16].

Another issue encountered in many distributed learning settings, and predominantly in popular FL methods, is the reliance on a central server, which effectively restricts such algorithms to an enterprise-only setting. Fully decentralised, or peer-to-peer algorithms, operate without a central authority and depend instead on the network topology [17, 18]. Nevertheless, in such cases enhanced robustness against Byzantine users is typically achieved only for dense network topologies, which leads to increased communication overhead [19]. This challenge is already evident for training or fine-tuning of large models within the FL framework, even in the absence of malicious players [20].

Ensemble methods combine the knowledge of multiple models, built with different architectures, parameters, and amount of available data, to solve a single task [21]. More precisely, ensembles are built on the strengths of every contributor in order to overcome the weaknesses of the individual models. These methods enlarge the representation space of a model, which is also one of the goals of transfer learning methods [22–25], notably deployed in foundation models [6]. Ensemble methods have been typically considered in a centralised setting, and thus implicitly assume that none of the participants are Byzantine. The problem of Byzantine par-

ticipants has recently been popularised through the widespread adoption of blockchains [26], where network participants vote on the validity of a set of transactions to reach agreement. This is implemented through a consensus protocol, which aims to provide both safety and liveness guarantees. Informally, this means that the system is able to deal with misbehaving participants.

In this work, we introduce a novel ML paradigm that combines consensus protocols with common ensemble learning methods, which we naturally term *consensus learning*. As opposed to the FL setup, individual participants do not share information about their ML models, nor any local data, which allows us to bypass data privacy and leakage issues. Instead, participants are required to only share their predictions for any given data inputs of a test dataset, to which any participant may contribute. To strengthen Byzantine robustness, these predictions then enter a communication phase, which is governed by a consensus protocol to ensure that the network reaches agreement. Here, honest participants will truthfully follow the rules of the consensus protocol, while malicious ones will attempt to stir the network to their desired outputs.

1.1 Main contributions

Consensus learning enhances typical ensemble weighting methods through a communication phase, where participants share their outputs until consensus is reached. We present a theoretical analysis of the performance of a consensus learning algorithm specialised to a binary classification task, which represents a toy model suitable for explaining the subtleties of this novel paradigm. More precisely, we deploy the Slush consensus protocol from the Snow family of protocols [27], which is a family of gossip protocols.

Our analysis provides lower bounds on the accuracy of a binary classifier deploying consensus learning and indicates which practical settings are suitable for using this type of algorithms. Moreover, we compare the Slush consensus learning algorithm with simple ensemble methods such as centralised majority rules and describe scenarios where the former can be the better performer. This analysis is supplemented by various numerical simulations, which describe the resilience of consensus learning against Byzantine participants.

While this work focuses mostly on classification tasks, consensus learning algorithms can also be applied to regression problems. There, robust local aggregation rules need to be deployed, similarly to Byzantine ML algorithms [9]. We also comment briefly on the applicability of consensus learning to unsupervised or self-supervised learning, and leave a detailed analysis of these use cases for future work.

1.2 Related works

Consensus learning is closely related to meta-learning [28], also known as learning-to-learn [24]. These methods involve a meta-learner, which is an ML model trained on the outputs of the base learners. In fact, weighting methods where the weights are based on the precision of the base learners can already be thought of as simple meta-learning methods. The first stage of consensus learning methods is, indeed, identical to that of meta-learning. The communication phase, however, does not involve a secondary training round per se; regardless, the analogy with meta-learning methods might stem from the local aggregation rules used in this phase, which can involve weighting methods.

Meta-learning methods have recently been used for unsupervised learning [29], indicating that consensus learning methods could be applied in such contexts as well. Of particular relevance would be unsupervised ensemble methods, such as consensus clustering [30], which offer a natural playground for extending consensus learning methods. Note that consensus

clustering algorithms combine clusterings from multiple sources without accessing the private data of individual participants. Nevertheless, this aggregation is done in a centralised fashion through a *consensus function*, which also assumes that participants are honest. A peer-to-peer adaptation can be implemented by adding a communication phase, where the consensus function would be used for local aggregations.

Another related approach to consensus learning is federated one-shot learning, which is an FL method that allows the central server to learn a model in a single communication round [31, 32]. Federated one-shot learning could be also classified as an ensemble method with some additional features, such as cross-validation user selection: only users achieving a baseline performance on some validation data can be part of the global model. Consensus protocols are generally not feasible for multi-round federated learning algorithms, due to the amount of computing resources required for obtaining satisfying results [19], but could be used for one or few-shot FL. Weaker forms of consensus, such as approximate agreement [33] and averaging agreement [19] have been argued to be more suitable for modern ML applications.

Knowledge distillation (KD) [34] based methods also share some similarities with consensus learning. There, the goal is to compress the knowledge of a group of *teacher models* into a smaller *student model* that can approximate the teachers’ predictions with high precision. In KD-based FL algorithms, participants only communicate their predictions on an unlabelled public test set, which are then used for improving local models [35] – see also *e.g.* [36, 37] and references therein for other recent works on the topic. Note that most research in this direction focuses on centralised algorithms.

Consensus algorithms have recently been deployed in distributed neural networks. In this setup, data is split among a number of agents, which share the same initialised neural networks. These models are then trained locally, with the local updates aggregated based on the network topology and a deterministic consensus algorithm. The setup is somewhat similar to peer-to-peer FL, and has been shown to achieve convergence using a small number of communication rounds after each training phase – see *e.g.* [38]. This setup becomes closer to our approach when using a heuristic adaptive consensus algorithm, which deploys local aggregation functions with varying weights [39]. This approach is based on switching communication graphs for the network topology, which is reminiscent of temporal graph neural networks [40].

Blockchain mechanisms, such as *proof-of-learning* (PoL) [41], also appear to share some of the features of consensus learning. The implementation described in [41], which was inspired by Kaggle machine learning competitions, involves a set of nodes called *trainers* that submit ML models for tasks previously set by other nodes, referred to as *suppliers*. These models are ranked according to their performance on unseen data by a set of randomly selected *validators*. A similar proposal was sketched out in [42]. Nevertheless, these PoL mechanisms still lack a basic utility: rather than aiming to collaboratively solve a problem, each node only tries to have the best model on their own. Blockchain empowered FL methods, such as those of [17, 43, 44], align more closely with our proposal; regardless, such blockchain implementations of FL methods are currently not feasible due to the shortcomings of EVM-based chains.¹ Additionally, certain client selection algorithms might be needed to ensure a minimum algorithm precision.

1.3 Organisation

The rest of the paper is organised as follows. In Section 2 we cover some background material on jury problems and ensemble learning, focusing on binary classification problems. We also review consensus mechanisms and their fundamental properties. Section 3 introduces the

¹EVM chains do not natively support floating point numbers.

principles of consensus learning and summarises the key results of our work. Section 4 provides a theoretical analysis for binary classifiers built with consensus learning algorithms, as well as a performance analysis in the presence of Byzantine learners. This analysis is extended in Section 5 through numerical simulations on non-iid² data. Finally, we summarise our findings and future research directions in Section 6.

2 Preliminaries

2.1 Ensemble learning

Ensemble methods provide powerful techniques for combining multiple ML models to make a common decision [45]. Algorithms developed using ensemble learning techniques are task-agnostic, thus generalising across a wide range of problems (see *e.g.* [46, 47]).

Definition 2.1 (Ensemble). *An ensemble is a collection of ML models whose predictions are combined together into a single output for any given input.*

Definition 2.2 (Base learner). *A base learner is one of the individual components of an ensemble.*

A base learner is thus a trained ML model deployed by a single participant of a network. The main premise of ensemble methods is that the errors of a single learner are compensated by the other learners [21]. In a distributed setting, such methods lead to significantly higher computational power and larger training datasets. Moreover, these techniques can reduce the risk of overfitting and increase the robustness of a model [48–50].

The effectiveness of ensemble learning techniques can also be appreciated through the perspective of hypothesis spaces [51]. Model building in supervised learning algorithms generally involves a search through a task-dependent hypothesis space, which can be understood as the set of functions between the input features and the output labels. In many instances, it is highly likely that the optimal hypothesis lies outside the hypothesis space of a single model. By combining multiple models, ensembles enlarge these individual spaces, thus increasing the likelihood of finding the optimal hypothesis [49].

This idea is closely related to the concept of *domain generalisation*, where multiple data sources are combined in order to improve a model’s generalisation performance on unseen target domains [52]. Additionally, ensemble learning techniques can achieve similar feats to *transfer learning* methods [23–25], whose goal is to leverage knowledge from different but related problems. The connection to these two concepts appears to be more pronounced especially when the base learners of an ensemble display significant dissimilarities, which will be a recurring theme of this work.

In ensemble learning, the outputs of the base learners can be combined together with two main methods: *weighting methods* and *meta-learning methods*. Typically, weighting methods turn out to be most suitable for cases when the performances of the learners are comparable. The simplest such method is a majority voting, where the assigned weights are all equal, and is commonly deployed in bootstrap aggregating (bagging) [53], or random forest [54] algorithms. More intricate weight assignments are used, for instance, in boosting algorithms [55].

Meta-learning algorithms, on the other hand, use a two-stage process in which the outputs of the base learners become inputs for an additional learning process [56]. Such methods are expected to perform extremely well when the base models have different performances on distinct subspaces of the dataset. Common meta-learning algorithms, such as *stacking* [56],

²Independent and identically distributed.

rely on a central server - the *meta-learner* - which is trained on the outputs of the base learners. Notably, both weighting methods and meta-learning algorithms are prone to various types of attacks from external users when considered in a distributed setting, which we aim to address in this work.

2.2 Jury problems

Ensemble methods implicitly assume that the base learners are *honest*. This assumption is partly relaxed in a decentralised setting, where a fraction of the base learners is allowed to be Byzantine.

Definition 2.3 (Honest participant). *An honest participant is a participant who follows the modelling process truthfully.*

Importantly, according to this definition, a low-performance model can still be labelled as honest. One of the fundamental results in ensemble methods is *Condorcet’s jury theorem*, which, despite its assumptions, captures essential aspects for building ML ensembles. This result is tailored to binary classification tasks, which we will also focus on for the rest of this work unless otherwise stated.

Hansen and Salamon [57] adapted the jury theorem to an ML context, by modelling a base learner as a Bernoulli trial X_i , with probability p_i of correctly identifying the label of a given input, for $i = 1, \dots, n$, where n is the number of base learners. In simplest terms, these success probabilities measure the accuracy of the base learners for the binary classification task at hand.

Definition 2.4 (Accuracy). *The accuracy of a classifier is the fraction of correctly identified samples in a test set.*

This measure can then be extrapolated to new inputs. Namely, the accuracy of an ML classifier will approximately give the probability of correctly identifying a new input. Of course, this assumes that the new input comes from the same distribution as the test data. In classification tasks, accuracy is based on the use of a *unit loss function*, which does not distinguish between false positives and false negatives. As such, accuracy may not be the best performance metric for imbalanced datasets, where metrics such as F1-score or area under the ROC curve provide better alternatives [58]. For our generic setting, however, accuracy will provide a reasonable metric.

To introduce Condorcet’s jury theorem, let us first define independence and homogeneity.

Definition 2.5 (Independence). *Base learners are independent if their associated random variables $\{X_i\}_{i=1, \dots, n}$, are pairwise independent.*

Definition 2.6 (Homogeneity). *A group of base learners is called homogeneous if all participants have the same accuracy on a specific input.*

The jury theorem dates back to the 18th century, and constitutes a majority rule ensemble method [57]: voters are given a binary choice and the collective decision corresponds to that of the majority. For simplicity, we assume that the number of voters is odd, to guarantee that a decision can always be made.

Condorcet’s jury theorem. *Given a homogeneous group of n independent base learners, for n odd, each having accuracy $p > \frac{1}{2}$, the accuracy \mathbb{P}_{Maj} of the ensemble built using a majority rule satisfies*

$$\mathbb{P}_{\text{Maj}}(p, n) = \sum_{j=\lceil \frac{n}{2} \rceil}^n \binom{n}{j} p^j (1-p)^{n-j} \geq p, \quad (2.1)$$

with equality for $n = 1$ or $p = 1$ only. Moreover, in the limit $n \rightarrow \infty$, we have

$$\mathbb{P}_{\text{Maj}}(p, n) \rightarrow 1 . \quad (2.2)$$

We refer to *e.g.* [59] for a proof of the first statement. The convergence in the large n limit follows from the law of large numbers and will be included in the proof of Proposition 4.2. The basic principle behind Condorcet’s jury theorem is that of the *wisdom of crowds*, *i.e.* the knowledge of a crowd is larger than that of a single member. This is, of course, not true in general, but convergence theorems can still be proven for heterogeneous or correlated juries – see *e.g.* [59, 60].

Despite their apparent simplicity, weighting methods are effectively used by many state-of-the-art ensemble methods, such as bagging [53] or boosting algorithms [61]. In fact, the simple majority rule used in Condorcet’s jury theorem turns out to be a very powerful aggregation rule for problems with a high degree of homogeneity, *i.e.* where the models of the base learners have similar performances. This was explicitly demonstrated by Nitzan and Paroush (Theorem 1 of [62]), who showed that the optimal (decisive) decision rule³ is a weighted majority, with the weights solely determined by the base learner accuracy.

This result shows, in particular, that in the homogeneous setting, the majority rule will outperform any other aggregation rule. Additionally, the majority rule will still perform close to optimal if the variance of the distribution of accuracies is not too large. In this sense, we introduce the notion of *diversity*, which we will revisit in subsequent sections.

Definition 2.7 (Diversity). *The diversity of a group of base learners is defined as the spread of the distribution⁴ of accuracies of the base learners.*

2.3 Consensus mechanisms

Consensus protocols were introduced in the context of distributed computing [63–65], to ensure a system’s security, resilience, and dependability [66], and today are at the heart of blockchains.

Blockchains are public databases which are updated and shared across many nodes in a network. Transactional data is stored in groups known as blocks, through the use of a unique identifier called a block hash, which is the output of a cryptographic hash function. This hash value is then part of the data of the next block of transactions, which thus links the blocks together in a chain. As such, each block cryptographically references its parent, and, thus, each block contains information about *all* previous blocks. As a result, the data in a block cannot change without changing all subsequent blocks, which would require the approval of the entire network. Every node in the network keeps a copy of the database and, consequently, every node must agree upon each new block and the current state of the chain. To accomplish this distributed agreement, blockchains use consensus mechanisms.

A protocol can solve the aforementioned problem of consensus if a set of conditions are satisfied [65–67]: every honest node eventually decides on some value (*termination*); if all nodes propose the same value, then all nodes decide on that value (*validity*); no node decides twice (*integrity*); no two honest nodes decide differently (*agreement*). These conditions become highly non-trivial in the presence of Byzantine nodes, and typically further assumptions about the environment might be needed for liveness and safety guarantees. These will be further discussed in Section 4.1.

³In the context of binary jury theorems, a decisive decision rule is a rule that ensures a decision for any set of juror votes.

⁴This is also known as *variability*.

Our primary focus will be on probabilistic consensus protocols, which rely on random or probabilistic processes. One such example is Nakamoto consensus [26] currently used by the Bitcoin network. Other examples include *gossip* protocols [68, 69], which are a class of communication protocols used to circulate information within a network. Probabilistic protocols typically require less communication overhead and thus tend to provide scaling advantages over deterministic variants; conversely, they can be subject to weaker resilience against malicious participants [70].

2.3.1 Snow consensus protocols

In the following, we consider the Snow family of consensus protocols [27]. These operate by repeatedly sampling the network at random, and steering correct nodes towards a common outcome, being examples of gossip protocols. In this section, we describe some of the technical aspects of the Slush protocol, which is the simplest protocol from this family. We include some new analytic results, summarised in Lemma 2.10 and Remark 2.9. Other technical details about the protocol are explained in Appendix A. See also [67, 71] for further recent analysis.

Consider a network consisting of n nodes and a binary query with the output choices formally labelled by two colours, red and blue. At the beginning of the Slush protocol, a node can either have one of the two coloured states or be in an uncoloured state. Then, each node samples the network at random, choosing k nodes to which they send a query. Every node responds to a query with its colour.

Once a node receives k responses, where typically $k \ll n$, it updates its colour if a threshold of votes $\alpha > \lfloor \frac{k}{2} \rfloor$ is reached. This process is repeated multiple times, and each node decides on the colour it ends up with after the last communication round. To ensure convergence, Slush needs $\mathcal{O}(n \log k)$ rounds [67], which is considerably lower than the $\mathcal{O}(n^2)$ rounds required in most deterministic protocols.

The dynamics of the Slush protocol can be modelled as a continuous-time Markov process [27]. For now, let us assume that all n participants are honest and that all nodes are in a coloured state. We will refer to \mathcal{S} as the state (or configuration) of the network at any given time. Without loss of generality, the state simply represents the number of blue nodes in the system and takes values in the set $\{0, \dots, n\}$. The process has two absorbing states, all-red and all-blue, corresponding to the final decision taken by the network.

Definition 2.8 (Slush absorption rates). *Let the absorption rates in the all-red (all-blue) state of the Slush protocol from the state with b blue nodes be \mathcal{R}_b (\mathcal{B}_b , respectively). These satisfy $\mathcal{R}_b + \mathcal{B}_b = 1$.⁵*

The exact expressions for these absorption probabilities are given in Appendix A, Corollary A.2. Based on these results, we make the following remark.

Remark 2.9. *The absorption probability in the all-blue state, \mathcal{B}_b , increases monotonically with the number of blue nodes, b .*

Additionally, it also follows that \mathcal{R}_b is a monotonically decreasing function with b . A new and important result that we will use throughout this paper follows. The proof can be found in Appendix A.

Lemma 2.10. *The Slush protocol is symmetric as long as all participants are honest, i.e.*

$$\mathcal{R}_b = \mathcal{B}_{n-b} , \tag{2.3}$$

⁵See also the remark at the end of Section 4.1.

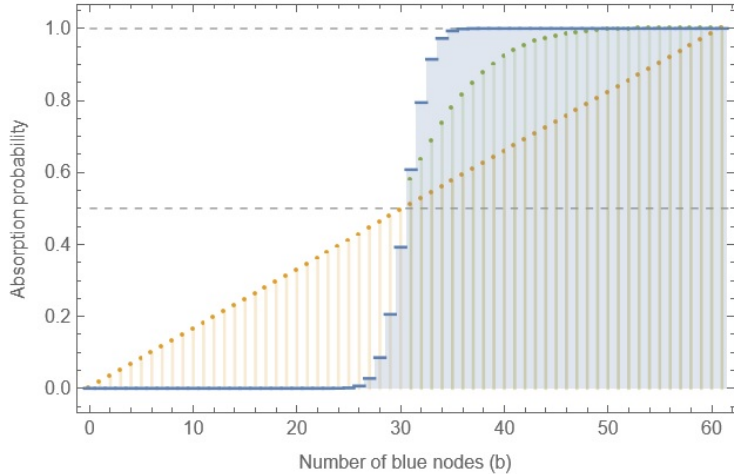


Figure 1: The \mathcal{B}_b absorption probability (in blue) for the Slush consensus protocol for $n = 61$, $k = 10$, $\alpha = 7$, as a function of the number of blue nodes b . In orange, the $\frac{b}{n}$ bound is plotted, while in green we have the Chvatal bound (4.8).

for any $b \in \{0, \dots, n\}$. Moreover, the majority absorption probability satisfies:

$$\mathcal{B}_b \geq \frac{b}{n}, \quad (2.4)$$

for $b \geq \lceil \frac{n}{2} \rceil$. Equality occurs for all b whenever $k = \alpha = 1$.

Figure 1 shows an explicit plot of the absorption probability \mathcal{B}_b and some relevant bounds for it. While the bound (2.4) is not particularly strong, it will play an important role in our analysis in Section 4. Note that this Lemma no longer holds in the presence of Byzantine nodes. Finally, let us mention that $\mathcal{R}_{b < \alpha} = 1$ (and $\mathcal{R}_{b > n - \alpha} = 0$), as, in these cases, the threshold for accepting a query can only be reached for the red (respectively blue) colour.

3 Consensus learning

In this section, we introduce consensus learning, a fully distributed ML paradigm that is based on consensus protocols. We focus on supervised ML methods, and briefly comment on how the algorithm could be adapted to unsupervised and self-supervised problems.

3.1 Algorithm description

Supervised consensus learning is a two-stage process that can be described as follows. In this description, we assume for now the existence of a global test set.

1. **Individual learning phase.** During the first stage, participants develop their own ML models, without the need to share any data or information about their models. At the end of this phase, participants determine their initial predictions for any given inputs.
2. **Communication phase.** During the communication phase, participants exchange their initial predictions on new inputs, and update them using a local aggregation function based on the outputs of the other base learners and their confidence in their own predictions. This phase is governed by a consensus protocol which may include several rounds, with the aim of guiding the network towards a common output. The outputs from the end of the communication phase will be the final outputs of the participants.

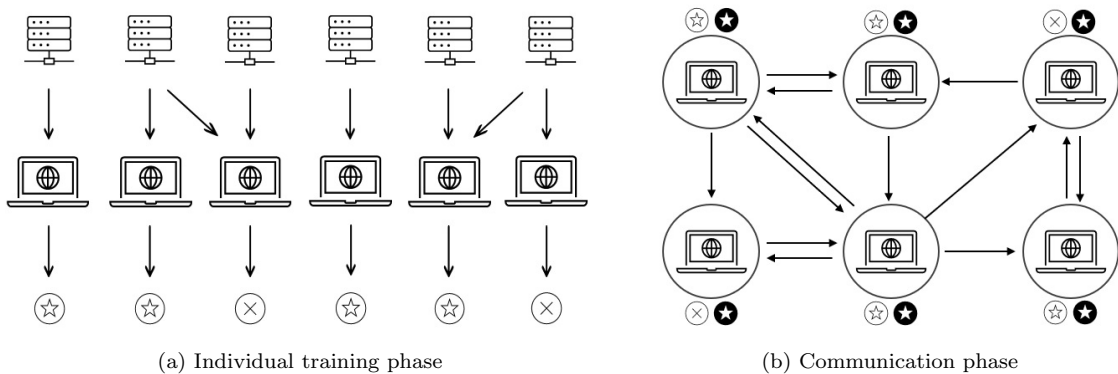


Figure 2: Supervised consensus learning. (a) In the the first stage, participants develop their own models, based on datasets that may overlap. At the end of this phase, each model determines an initial prediction (hollow circles) for any new input. (b) In the communication phase, the initial outputs are exchanged between the participants, which eventually reach consensus on a single output (filled circles).

The two stages are depicted in Figure 2. In an ideal case, the (honest) base learners reach consensus on a single global output in the communication phase. This unique output would then be the output of the ensemble formed by the base learners. However, this is not a requirement in fully decentralised algorithms, where participants may reach different final decisions.

Consensus learning also allows for direct implementations on decentralised platforms such as blockchains. There, participants may either propose data for testing using proof-of-stake based protocols or use a predefined test set provided by an independent party. The latter can be facilitated through the use of a smart contract on Ethereum Virtual Machine (EVM) compatible blockchains.

The algorithm can also be adapted to self-supervised or unsupervised ML problems, where participants only have access to (partly) unlabelled data. For the former, each participant would deploy self-supervised learning techniques during the individual learning phase, such as contrastive learning [72] or auto-associative learning [73]. Then, the communication phase would proceed similarly to the supervised setting; here, the test set could include data inputs from the training sets of the individual participants, with the ensemble outputs being used to improve local models.

Unsupervised ensemble methods, such as consensus clustering [30], combine clusterings for multiple sources without accessing the private data of individual participants. This aggregation is commonly done in a centralised fashion through a *consensus function*, which also assumes that participants are honest. Such methods can be adapted to peer-to-peer settings through the implementation of a communication phase, as described above. There, the consensus function would be used for local aggregations.

3.2 Summary of key results

Consensus learning methods are fully distributed ensemble techniques, being thus effective for generalising across a wide range of problems, as described in Section 2.1. Moreover, consensus learning is a much simpler type of meta-learning, being less demanding from a computational perspective. It is worth pointing out that meta-learning algorithms create additional sources of overfitting [74], which are absent in consensus learning algorithms.

Another advantage of consensus learning is that it preserves the explainability of ensemble weighting methods. Thus, the relatively simple design of such algorithms offers transparency and interpretability. Importantly, consensus learning does not rely on a single central server.

Additionally, adequate choices of probabilistic consensus protocols can result in low communication overhead, on par with that of centralised weighting methods. It is also natural to consider whether or not consensus-based methods can improve the performance of classical weighting methods, which we will discuss further below.

For classification tasks in supervised learning, perhaps the simplest weighting method is the equal-weight majority rule, deployed in methods such as Random Forests [54]. In the next section, we consider binary classification and present a theoretical analysis of the performance of a consensus learning algorithm against centralised majority rules. For this, we deploy the Slush consensus protocol [27] in the communication phase. We use accuracy as a performance metric (as per Definition 2.4) and analyse three types of scenarios, as discussed below.

I. Homogeneous scenario. Arguably the simplest scenario to consider from an analytical point of view is one where all base learners have the same accuracy. In this homogeneous setting, Condorcet’s jury theorem has long been one of the main pillars of ensemble learning. One of our most important results in this direction is a generalisation of Condorcet’s jury theorem to a consensus learning algorithm.

Theorem 1. *Consider a homogeneous group of n independent base learners, with accuracies p for a binary classification task. Then, the accuracy \mathbb{P}_S of the consensus learning algorithm using the Slush protocol satisfies:*

$$\mathbb{P}_S \geq p, \tag{3.1}$$

for any $p > \frac{1}{2}$, with equality only occurring for $n = 1$ or $p = 1$. Moreover, \mathbb{P}_S can be brought arbitrarily close to 1 for any $p > \frac{1}{2} + \frac{1}{n}$, and large enough n .

The proof of this statement is rather involved, and is left to Appendix B. Nevertheless, the proof only uses simple features that are specific to the Slush protocol and could thus be generalised to other probabilistic consensus protocols. Other main results in the homogeneous setting are discussed in Section 4.2 and include lower bounds on the accuracy of the consensus learning algorithm using the Slush protocol, as well as comparisons with majority and supermajority rules.⁶

II. Partly heterogeneous scenario. Adaptations of jury theorems to heterogeneous juries have previously been discussed in the literature – see *e.g.* Theorem 4 of [60], Theorem 3 of [75], as well as Theorem 3 of [76]. In fact, we expect that a similar result to Theorem 1 would hold for heterogeneous groups, but we do not explicitly pursue this direction.

Instead, we compare consensus learning algorithms with majority rules in partly heterogeneous settings: more precisely, the accuracies of the base learners will be split into distinct homogeneous groups. We will see, in particular, that consensus learning algorithms can perform better than majority rules as long as the learners are *diverse* enough, *i.e.* when the distribution of accuracies of the base learners has a certain degree of heterogeneity, as per Definition 2.7. This topic will be discussed in more detail in Section 4.3.

III. Almost homogeneous scenario with Byzantine nodes. An implicit assumption of well-established ensemble methods is that the base learners are honest. In a fully decentralised setting (or peer-to-peer), this assumption is relaxed, due to the presence of Byzantine participants. We clarify this notion in our framework below.

⁶Supermajority rules are essential for proving Theorem 1.

Definition 3.1 (Byzantine participant). *A malicious or Byzantine participant is a participant who is not honest. This also extends to the communication phase of a consensus learning algorithm.*

A Byzantine participant may share any outputs they wish during the communication phase. These will usually be decided based on their adversarial strategy. Section 4.4 will present some analytical results regarding consensus learning algorithms with Byzantine users, where the honest base learners will form a homogeneous group. Our analysis comprises a comparison with majority and supermajority rules in the presence of Byzantine users.

Numerical results. To better understand these scenarios and to provide more supporting evidence in favour of consensus learning algorithms, we will also present various numerical simulations in Section 5. We will use non-iid data from the LEAF benchmark [77] and will analyse the effect of Byzantine users in more detail. Furthermore, we present slight modifications of the Slush protocol that can improve algorithm performance and Byzantine resilience.

4 Theoretical analysis

In this section we present a theoretical analysis of a consensus learning algorithm using the Slush consensus protocol, specialised to a binary classification problem. The underlying assumptions of this analysis are presented in Section 4.1. For ease of notation, we introduce the following terminology.

Definition 4.1 (Slush algorithm). *The Slush algorithm is the consensus learning method deploying the Slush consensus protocol in the communication phase.*

4.1 Modelling premises

To provide an in-depth analysis of a typical consensus learning algorithm, we will make some simplifying assumptions. Our analysis will be an extension of the Hansen and Salamon adaptation of jury theorems to an ML context [57]. As such, we will model base learners as Bernoulli trials; the success probabilities of these random variables correspond to a chosen performance metric in the binary classification task. Another modelling assumption will be the independence of the base learners, as per Definition 2.5. We will relax this assumption in Section 5, where we conduct a simulation of a numerical consensus learning algorithm. Let us also mention that the number of base learners will typically be assumed odd unless otherwise stated.

For a concrete illustration of the communication phase, we will consider the Slush consensus protocol [27], briefly summarised in Section 2.3.1. We also refer to Appendix A for a more technical discussion of this protocol. This consensus protocol can be adapted to a binary classification task as follows. Consider a query, representing a data entry which needs to be classified by the ensemble. Initially, each node picks a class (which we refer to as *colour*) as dictated by their local ML algorithm. For the binary classification problem of interest, we assume without loss of generality that class 1 (labelled as *blue* colour) is the correct class of some given input to be classified, as opposed to class 0 (labelled as *red* colour). Thus, at the beginning of the protocol, each node will be in a coloured state. These states can then change during the communication rounds.

Consensus protocols assume a form of *synchronicity*, as agreement cannot be reached in a fully asynchronous setting [78]. For the technical analysis of the consensus learning algorithms,

we will also assume a synchronous setting in the communication phase. Additionally, we will assume that consensus is eventually reached within the network, such that all (honest) base learners decide on the same final output.

Finally, let us briefly comment on the presence of Byzantine nodes. As pointed out in [27], if the number of Byzantine nodes is greater than the threshold α of accepting a query, then the Markov chain modelling the Slush protocol appears to have only a single absorbing state. Of course, in practical terms, the transition probabilities away from the all-blue state would be arbitrarily small. Nevertheless, throughout the paper, we will not consider this case. Another subtlety concerns the identity $\mathcal{R}_b + \mathcal{B}_b = 1$ for the absorption probabilities. It was recently argued that adversarial strategies that assume knowledge of the whole network can, in fact, indefinitely stall the protocol [79]. Our analysis will be limited to fixed (extreme) adversarial strategies, where the Byzantine nodes will communicate the wrong class at all times. In such cases, the protocol can still be modelled as a Markov chain with fixed transition rates, thus ensuring that there is no closed communicating class, apart from the two absorbing states.

4.2 Homogeneous case

As a first scenario, we consider the homogeneous case, where each participant j has the same accuracy in classifying a new input, $p_j = p$, for $j \in \{1, \dots, n\}$, with $p \in [0, 1]$. We will discuss the odd $n \in \mathbb{N}$ case below.

4.2.1 Majority rules

In the homogeneous context, the first natural question is how a consensus learning algorithm will compare to a single implementation of a majority rule. Theorem 2 gives a first result in this direction.

Theorem 2. *Given a homogeneous group of n independent base learners, each with accuracy p for a binary classification problem, the majority rule will outperform⁷ the Slush algorithm, as long as $p > \frac{1}{2}$ and $\alpha \neq \lceil \frac{n}{2} \rceil$. The Slush algorithm will achieve the same accuracy only for $\alpha = \lceil \frac{n}{2} \rceil$.*

The veracity of this affirmation can be inferred from the Nitzan-Paroush theorem on optimal decision rules (*i.e.* Theorem 1 of [62]). This states that the majority rule is the optimal decisive decision rule in the homogeneous setting. However, it is not entirely clear whether that theorem holds for our algorithm due to the existence of a communication phase. Thus, we provide below an alternative proof for this theorem.

Proof. Note first that the number of blue nodes b before the communication phase starts follows a binomial distribution. Thus, the probability of success for the Slush protocol is given by:

$$\begin{aligned} \mathbb{P}_S(n, k, \alpha, p) &= \sum_{b=0}^n \binom{n}{b} \mathcal{B}_b p^b (1-p)^{n-b} \\ &= \sum_{b=\alpha}^{n-\alpha} \binom{n}{b} \mathcal{B}_b p^b (1-p)^{n-b} + \sum_{b=n-\alpha+1}^n \binom{n}{b} p^b (1-p)^{n-b}. \end{aligned} \tag{4.1}$$

We would like to compare this expression with the expression (2.1) for the homogeneous majority rule. We immediately see that for $\alpha > \frac{n}{2}$, the first sum in (4.1) vanishes, and we

⁷That is, the accuracy of the majority rule is larger than that of the Slush algorithm.

have:

$$k \geq \alpha \geq \left\lceil \frac{n}{2} \right\rceil : \quad \mathbb{P}_S(n, k, \alpha, p) \leq \mathbb{P}_{\text{Maj}}(p, n) , \quad (4.2)$$

with equality only for $\alpha = \left\lceil \frac{n}{2} \right\rceil$. The more interesting case to analyse is $\alpha < \frac{n}{2}$. For this, note that the first sum in the bottom line of (4.1) can be further decomposed into

$$\sum_{b=\alpha}^{n-\alpha} \binom{n}{b} \mathcal{B}_b p^b (1-p)^{n-b} = \sum_{b=\alpha}^{\left\lfloor \frac{n}{2} \right\rfloor} \binom{n}{b} \mathcal{B}_b p^b (1-p)^{n-b} + \sum_{b=\left\lceil \frac{n}{2} \right\rceil}^{n-\alpha} \binom{n}{b} \mathcal{B}_b p^b (1-p)^{n-b} . \quad (4.3)$$

Then, defining

$$\Delta \mathbb{P} = \mathbb{P}_{\text{Maj}}(p, n) - \mathbb{P}_S(n, k, \alpha, p) , \quad (4.4)$$

we find that for $n = 2m + 1$, $\Delta \mathbb{P}$ reduces to:

$$\Delta \mathbb{P} = \sum_{b=m+1}^{n-\alpha} \binom{n}{b} \left(\mathcal{R}_b p^b (1-p)^{n-b} - \mathcal{B}_{n-b} p^{n-b} (1-p)^b \right) . \quad (4.5)$$

To make further progress, we consider the ratio of the individual terms in the above summation:

$$\kappa_b \equiv \frac{\mathcal{R}_b}{\mathcal{B}_{n-b}} \times \left(\frac{p}{1-p} \right)^{2b-n} , \quad (4.6)$$

for $m + 1 \leq b \leq n - \alpha$. Using Lemma 2.10, and since $2b > n$ for the range of interest, we have $\kappa_b > 1$ as long as $p > \frac{1}{2}$, with $\kappa_b = 1$ at $p = \frac{1}{2}$ and $\kappa_b < 1$ otherwise. This concludes our proof. \square

While the Slush algorithm cannot improve on the accuracy of an already optimal decisive decision rule in the homogeneous setting, we would still like to find a lower bound on its accuracy to illustrate its functionality. We thus seek a generalisation of Condorcet's jury theorem, which compares the accuracy of the ensemble generated through the Slush algorithm with that of a single base learner. The first part of this statement, illustrated by Theorem 1, already gives such a lower bound, namely $\mathbb{P}_S \geq p$, for $p > \frac{1}{2}$. In simple terms, this states that the ensemble built using this fully decentralised paradigm is more accurate than any of the base learners. A different bound on the Slush accuracy can be determined as follows.

Proposition 4.2. *Consider a homogeneous group of n independent base learners, each with accuracy p . For large enough n , a lower bound for the accuracy of the Slush algorithm is given by*

$$\mathbb{P}_S(n, k, \alpha, p) > 1 - e^{-2\left(\frac{\alpha}{k} - \frac{1}{2}\right)^2 k} , \quad (4.7)$$

as long as $p > \frac{1}{2}$.

Proof. The proof of the statement uses the Chvatal tail bounds of the hypergeometric distribution [80]. In particular, Theorem 1 of [27] shows that

$$\mathcal{R}_b \leq e^{-2\left(\frac{\alpha}{k} - 1 + \frac{b}{n}\right)^2 k} , \quad (4.8)$$

where $b \geq \left\lceil \frac{n}{2} \right\rceil$. Then, we have

$$\mathbb{P}_S > \sum_{b=\left\lceil \frac{n}{2} \right\rceil}^n \mathbb{P}(S_n = b) (1 - \mathcal{R}_b) > \sum_{b=\left\lceil \frac{n}{2} \right\rceil}^n \mathbb{P}(S_n = b) \left(1 - \mathcal{R}_{\left\lceil \frac{n}{2} \right\rceil} \right) , \quad (4.9)$$

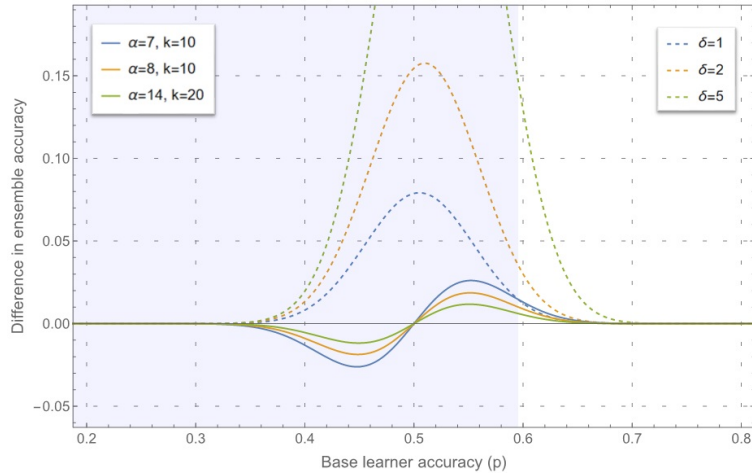


Figure 3: Solid lines: difference in accuracy between the simple majority rule and the homogeneous Slush algorithm, $\Delta\mathbb{P}$, with $n = 101$, against the base learner accuracy p . Dashed lines: differences in accuracy between majority and δ -supermajority rules. The shaded area shows the region where the Slush algorithm with $\alpha = 7$, $k = 10$ outperforms the $\delta = 1$ supermajority.

with $S_n = \sum_{i=1}^n X_i$ being the sum of the Bernoulli trials associated to the base learners. For the last inequality, we use the fact that the value $b = \lceil \frac{n}{2} \rceil$ minimizes the expression $(1 - \mathcal{R}_b)$ for the range of the sum. Furthermore, it is not difficult to see that

$$\mathbb{P}\left(S_n \geq \lceil \frac{n}{2} \rceil\right) = \mathbb{P}\left(\frac{S_n}{n} > \frac{1}{2}\right) \geq \mathbb{P}\left(\left|\frac{S_n}{n} - p\right| < p - \frac{1}{2}\right), \quad (4.10)$$

which, by the weak law of large numbers, converges to 1 for large n . The result follows immediately. \square

Tail bounds are rather conservative, and thus the actual accuracy of the Slush algorithm is expected to be considerably better than this bound. Nonetheless, the bound provides a different perspective on the performance on the Slush algorithm. Specifically, for accurate classifiers with $p > \frac{1}{2}$, this lower bound improves as the threshold parameter α for accepting a query increases. On the other hand, it was argued in [67] that values of α closer to $k/2$ are more suitable for Byzantine consensus protocols. This is thus an important distinction between the objectives of consensus protocols in distributed computing and those of protocols designed for consensus ML. Nevertheless, Theorem 1 shows that consensus learning can leverage protocols primarily designed to safeguard distributed networks.

Other finite n bounds, similar in spirit to that of Proposition 4.2, can be found using bounds for the binomial distribution [81]. However, we would like to find a stronger result for the case of large n . Ultimately, this search concludes with the second statement of Theorem 1, which conveys that the Slush algorithm is indeed an efficient algorithm. The proof of the statement relies on supermajority rules, which we introduce next.

4.2.2 Supermajority rules

The simple majority rule discussed so far requires that more than half of the votes are cast for an output. Supermajority rules increase this acceptance threshold and can lead to enhanced stability in the voting process, as well as increased legitimacy in the final decision [82].

Definition 4.3 (δ -supermajority rule). *A δ -supermajority rule is a majority rule for which the acceptance threshold required to choose an alternative is $\lceil \frac{n}{2} \rceil + \delta$. Alternatively, one can use a fraction q of the votes, with $\lceil qn \rceil$ votes required for taking a decision.*⁸

The following theorem shows that the Slush algorithm can outperform any $\delta \geq 1$ supermajority rule, even in a homogeneous setting. This is rather noteworthy since low δ supermajority rules are still very close to being optimal decision rules.

Theorem 3. *For any $\delta > 0$, there exists a value $p_{\text{th}}(\delta) > \frac{1}{2}$ such that the accuracy of Slush algorithm built with a homogeneous group of independent base learners with accuracies $p \leq p_{\text{th}}(\delta)$ will be larger than the accuracy of a δ -supermajority rule.*

The proof of this theorem can be found in Appendix B. To get a grasp on how the threshold value p_{th} changes with δ , we perform a numerical analysis below.

Numerical bounds. Figure 3 shows a comparison between majority aggregation rules and the Slush algorithm, for a homogeneous group of independent learners. We highlight that the Slush algorithm appears to perform exceptionally well compared to supermajority rules, even for $\delta = 1$. This indicates that the bound found in Theorem 3 should increase rapidly with δ . Table 1 gives some numerical values⁹ for the threshold value p_{th} , approximated to two decimal places, for $n = 101$ and a varying δ . These values validate our expectations that p_{th} increases rather fast with δ .

δ	0	1	2	3	4
{6, 10}	0.5	0.56	0.62	0.68	0.73
{7, 10}	0.5	0.6	0.68	0.76	0.83
{8, 10}	0.5	0.63	0.74	0.83	> 0.87
{14, 20}	0.5	0.7	0.84	0.87	> 0.88

Table 1: Lower bounds on the threshold values $p_{\text{th}}(\delta)$ for $n = 101$ and varying values of $\{\alpha, k\}$, as indicated in the first column.

Remarkably, supermajority rules can still be shown to satisfy jury theorems, as long as the base learner accuracy is larger than the acceptance threshold – see Theorem 2 of [76]. These results allow us to improve on the “large n ” behaviour of the Slush algorithm from Proposition 4.2. In the following, we shall use the fraction $q > \frac{1}{2}$ of votes required to accept a proposal when discussing a supermajority rule, as introduced in Definition 4.3.

Lemma 4.4. *For large enough n , the accuracy of the Slush algorithm with homogeneous and independent base learners can be brought arbitrarily close to 1, if the base learner accuracy p satisfies*

$$q < p < p_{\text{th}}(q) , \tag{4.11}$$

for some $q > \frac{1}{2}$, where $p_{\text{th}}(q)$ is the value below which the Slush algorithm outperforms the q -supermajority rule.

This result combines Theorem 3 with Theorem 2 of [76], and thus the proof is straightforward. The latter claims that the q -supermajority rule leads to unit success probability for large n as long as $p > q$. To get a sense of how the interval in Lemma 4.4 evolves with n , we list approximate values of $p_{\text{th}}(q)$ in Table 2. Of course, it is not obvious that $p_{\text{th}}(q)$ remains

⁸This decision rule is not decisive, in the sense of footnote 3.

⁹More precisely, the values shown in Table 1 are lower bounds for p_{th} .

n	51	101	201	501
$p_{th}(q)$	0.92	0.88	0.8	0.73

Table 2: Lower bound on threshold value $p_{th}(q)$, for $q = 0.55$, $k = 10$, $\alpha = 7$ and varying n .

greater than q as n increases further. Resolving this issue is crucial for proving Theorem 1.

Lemma 4.4 makes a clear statement on the accuracy of the Slush algorithm for $q < p < p_{th}(q)$. We would like to extend this interval further, for any $p > p_{th}(q)$ values. For this, we use the following result.

Lemma 4.5. *The accuracy of the Slush algorithm with homogeneous and independent base learners is a strictly monotonically increasing function of the base learner accuracy.*

Proof. Let us define the function $F(b^*, n; p)$ as:

$$F(b^*, n; p) = \sum_{b=b^*}^n \binom{n}{b} p^b (1-p)^{n-b} = 1 - \sum_{b=0}^{b^*-1} \binom{n}{b} p^b (1-p)^{n-b}. \quad (4.12)$$

The Slush algorithm accuracy defined in (4.1) can be also expressed in terms of this function as follows:

$$\mathbb{P}_S(p) = \mathcal{B}_0 F(0, n; p) + (\mathcal{B}_1 - \mathcal{B}_0) F(1, n; p) + \dots + (\mathcal{B}_n - \mathcal{B}_{n-1}) F(n, n; p). \quad (4.13)$$

The veracity of this statement can be checked backwards: from (4.13) we can collect the terms $p^b (1-p)^{n-b}$ for all b and see that (4.1) is recovered. Note that the coefficients of all $F(b^*, n; p)$ terms are positive due to Remark 2.9. As such, if $F(b^*, n; p)$ were an increasing function of p , then so would $\mathbb{P}_S(p)$. To show this, we look at the first derivative

$$\begin{aligned} \frac{d}{dp} F(b^*, n; p) &= - \sum_{b=0}^{b^*-1} \binom{n}{b} p^{b-1} (1-p)^{n-b-1} (b-np) \\ &= \binom{n}{b^*} b^* p^{b^*-1} (1-p)^{n-b^*} > 0, \end{aligned} \quad (4.14)$$

which is clearly positive. The proof of the identity used on the second line can be found in Appendix A, Lemma A.3. This concludes the proof. \square

Towards proving Theorem 1. Equipped with the above results, we are now ready to sketch the proof for Theorem 1. The full proof of the theorem can be found in Appendix B.

Sketch of Proof. The first statement of the theorem ($\mathbb{P}_S \geq p$) follows rather simply from an application of Lemma 2.10. Nevertheless, a proof of the second statement ($\mathbb{P}_S \rightarrow 1$) requires multiple ingredients. This statement builds on Theorem 3, according to which large accuracies for the Slush algorithm can occur whenever the base learner accuracy p is in the interval $(q, p_{th}(q))$. Using the monotonicity of the Slush algorithm proved in Lemma 4.5, we can eliminate the upper bound of this interval.

However, the only remaining issue is showing that $q < p_{th}(q)$, for any n , such that Theorem 3 is valid. This part of the proof is rather tedious, and can be found in Appendix B. \square

4.3 Diversifying the base learners

The Slush algorithm combines the distribution of outputs of the base learners with an additional random variable responsible for the communication phase. For a general heterogeneous setting, the binomial distribution from the homogeneous problem is replaced by a Poisson binomial distribution. The ratio κ_b introduced in the proof of Theorem 2 was shown to dictate the behaviour of the Slush algorithm, as compared to the majority rule. For future reference, we define it more formally below.

Definition 4.6 (Control ratio). *Consider a learning problem with n base learners. Let b be the number of blue nodes at the start of the communication phase, and let $\mathbb{P}(S_n = b)$ be the probability that this state can arise from the initial outputs of the base learners, where $S_n = \sum X_i$ is the sum over the Bernoulli random variables assigned to the base learners. Then, we define the control ratio as:*

$$\kappa_b = \frac{\mathcal{R}_b}{\mathcal{B}_{n-b}} \times \frac{\mathbb{P}(S_n = b)}{\mathbb{P}(S_n = n - b)}, \quad (4.15)$$

where \mathcal{R}_b (respectively \mathcal{B}_b) are the absorption probabilities in the all-red (all-blue, respectively) states, starting from a state with b blue nodes.

This particular definition of the control ratio will be relevant for values $b \geq \lceil \frac{n}{2} \rceil$. The control ratio can be used to deduce a set of simple sufficient conditions for the Slush algorithm to outperform the majority rule: $\kappa_b < 1$. The argument for this statement is very similar to that presented in the proof of Theorem 2. As such, a simple analysis of this ratio leads to insights into the type of problems where the Slush algorithm would be more suitable. Schematically, we can interpret the control ratio as

$$\kappa = \left(\begin{array}{c} \text{Asymmetry of} \\ \text{consensus protocol} \end{array} \right) \times \left(\begin{array}{c} \text{Diversity of} \\ \text{base learners} \end{array} \right). \quad (4.16)$$

Here, the *diversity* is measured by the variance of the distribution of accuracies, as per Definition 2.7. Based on these two factors, we have the following two cases of interest:

- Asymmetric problems, *i.e.* problems which break the symmetry of the Slush protocol from Lemma 2.10 in favour of the correct output.
- Heterogeneous problems, *i.e.* situations in which we deal with a diverse group of learners, which include both strong and weak learners.

In view of these conclusions, we consider a semi-homogeneous setting, with two performance groups. The main result of this subsection is the following theorem.

Theorem 4. (Performance groups) *Consider a binary classification task with two homogeneous groups of classifiers of sizes n_1 and n_2 with $n_1 > n_2$, having accuracies p_1 and p_2 , respectively. Assume that the classifiers are independent and let $n = n_1 + n_2$, with n odd. Then, for $p_2 < \frac{1}{2}$, there exists $p_{\max} > \frac{1}{2}$ such that the Slush algorithm outperforms the majority rule for any p_1 in the region $[0, p_{\max}]$, where p_{\max} is bounded by*

$$p_{\max} \leq \frac{\lceil \frac{n}{2} \rceil}{1 + n_1}, \quad (4.17)$$

with the upper bound reached for the limiting case $p_2 = 0$.

The proof of this statement is left to Appendix B. Let us stress that in this setting the weak learners will still truthfully follow the modelling process. Thus, they can change their state in the communication phase, according to the majority of the sampled nodes, even if $p_2 = 0$. In Section 5, we will generalise this setting to a completely heterogeneous one. There, we will see that the Slush algorithm can significantly outperform a centralised majority rule.

4.4 Byzantine tolerance in consensus learning

Byzantine fault-tolerant consensus protocols are known for their ability to tolerate failures and withstand malicious attacks on a network. Accordingly, the next logical problem deserving our attention is that of a network which includes such Byzantine participants, as introduced in Definition 3.1.

Malicious nodes will usually follow an *adversarial strategy*, with the aim of steering the network towards their preferred outcome. Note that this is different from the scenario described by Theorem 4, as, in that case, all participants behave the same during the communication phase. To describe the behaviour of the network in the presence of such participants, we can consider the extreme scenario in which the Byzantine participants know the correct outcome with certainty, but decide against sharing it.

Definition 4.7 (Perfectly Byzantine participant). *A perfectly malicious (or perfectly Byzantine) participant is a participant who knows with certainty the correct label of any data inputs, and who, when queried, will always respond with the wrong label.*

Given this definition, we can now consider a learning problem which includes a number of perfectly malicious participants. Let us point out that the scenario described by the following theorem is still one with a high degree of homogeneity; accordingly, the majority rule is expected to perform rather well.

Theorem 5. (*Perfectly malicious nodes*) *Consider a group of $f < \alpha$ perfectly malicious participants, in an otherwise homogeneous group of independent base learners of size n . Let the accuracy of the $c = n - f$ honest base learners be p . Then, the majority rule will outperform the Slush algorithm with parameters k and α , as long as $p > \frac{1}{2}$ and $\alpha \neq \lceil \frac{n}{2} \rceil$. The Slush algorithm will achieve the same accuracy only for $\alpha = \lceil \frac{n}{2} \rceil$.*

This theorem is the analogous of Theorem 2 to the scenario involving Byzantine nodes. Its proof makes use of the exact form of the absorption probabilities in the Slush protocol in the presence of Byzantine nodes, as described by Lemma A.5, which can be found in Appendix A.

The observant reader may have noticed the constraint $\alpha > f$ used in the previous theorem. Recall that we model the Slush protocol as a continuous-time Markov chain, with two absorbing states. The death and birth rates are given by the probability that the next query changes a node's colour for the red or blue colours, respectively. If f were to be larger than or equal to α , then the state with all honest nodes being blue would no longer be an absorbing state for this process. Thus, to avoid this complication, we set $\alpha > f$ in the statement of the above theorem. Note, however, that in a practical scenario, the protocol will only run a finite amount of rounds. Hence, even in the $\alpha \leq f$ case, there would be a finite probability of reaching consensus for the blue colour. This can be modelled by artificially setting the death rate from the all-blue state to zero.

Proof of Theorem 5. The majority rule and Slush algorithm accuracies are given by:

$$\begin{aligned} \mathbb{P}_{\text{Maj}}(n, c, p) &= \sum_{b=\lceil \frac{n}{2} \rceil}^c \binom{c}{b} p^b (1-p)^{c-b}, \\ \mathbb{P}_{\text{S}}(n, c, k, \alpha, p) &= \sum_{b=0}^c \binom{c}{b} \mathcal{B}_b p^b (1-p)^{c-b}. \end{aligned} \tag{4.18}$$

From here, we again look at $\Delta\mathbb{P} = \mathbb{P}_{\text{Maj}} - \mathbb{P}_{\text{S}}$, and apply the same methods as in the proof of Theorem 2. As before, if $\alpha > \lceil \frac{n}{2} \rceil$, the majority rule outperforms Slush, while the performances are identical for $\alpha = \lceil \frac{n}{2} \rceil$. Then, using Lemma A.5, it follows that the control ratio

satisfies:

$$\kappa_b > \frac{n-b}{c-b} \times \left(\frac{p}{1-p} \right)^{2b-n} > 1, \quad (4.19)$$

for $\lceil \frac{n}{2} \rceil \leq b \leq n - \alpha$.

□

Faulty communication. Theorem 5 offers some insight into how the homogeneous Slush algorithm performs in the presence of Byzantine nodes. In a more realistic scenario, Byzantine behaviour can simply be due to faulty communication. We can expect, in particular, that such participants will only be able to send at most one query, but will not be able to receive any responses. As a result, such participants will effectively drop out in the communication phase. This shows another advantage of consensus learning over centralised ensemble methods since the former can identify this type of faulty participants.

Definition 4.8 (Faulty participant). *A faulty participant is a participant who does not participate in the communication phase.*

Definition 4.9 (Perfectly faulty participant). *A perfectly faulty participant is a faulty participant whose initial output is incorrect.*

We remind the reader that a perfectly faulty participant is not trying to stall the system. Instead, their ML model can be thought of as having very low accuracy (zero) – see also Section 4.1. The communication phase of the Slush algorithm will be able to detect the faulty participants, which will thus drop out. However, the initial responses of the faulty nodes should still be considered in a majority rule, as there would be no means of identifying a faulty connection in such cases. Thus, introducing perfectly faulty participants is equivalent to considering supermajority rules instead of a simple majority.

Conjecture 4.10. *Consider a homogeneous group of c independent base learners, each with accuracy p . Let there be f perfectly malicious participants and f' perfectly faulty participants. For an appropriate choice of $f' > 0$, there exists a value $p_{\text{th}} > \frac{1}{2}$ such that the Slush algorithm with parameters $\alpha > f$ and $p \leq p_{\text{th}}$ outperforms the majority rule.*

As previously alluded to, it is not difficult to see that this statement reduces to a comparison between the Slush algorithm and δ -supermajority rules, with $\delta = f'$. As such, this proposition is an interpolation between Theorem 5 and Theorem 3. We do not have a proof of this statement, but we offer some numerical evidence below.

Remarkably, small values of $\delta = f'$ already ensure that the threshold value p_{th} is larger than 50%. Table 3 gives some numerical values for $p_{\text{th}}(\delta)$, approximated to two decimal places, for $n = 101$, and fixed protocol parameters.

δ	0	1	2	3
$f = 0$	0.5	0.7	0.84	0.87
$f = 1$	0.49	0.69	0.84	0.87
$f = 5$	0.47	0.67	0.83	0.87
$f = 10$	0.47	0.67	0.83	0.87

Table 3: Lower bounds on threshold value $p_{\text{th}}(\delta)$ for fixed $n = 101$, $k = 20$, $\alpha = 14$ and varying f and δ .

5 Numerical simulations

In this section, we present numerical simulations of the Slush consensus learning algorithm that extend beyond our previous analysis. It is worth pointing out that the communication phase in the previously described scenarios did not make use of the conviction of the base learners in their prediction. In this regard, we will also present a modified Slush algorithm which uses local parameters instead of globally defined ones. Such modifications can stir the network in favour of the better classifiers.

Definition 5.1 (Strong confidence). *A participant whose local model has accuracy $p > \frac{1}{2}$ is said to have strong confidence in their result if their local threshold parameter for accepting a query satisfies $\alpha \geq kp$.*

To simulate the behaviour of a realistic fully decentralised network, we will consider non-iid data. It should be noted that heterogeneity among the training sets has been shown to pose serious challenges for FL algorithms to achieve high levels of precision [83]. In Section 5.1 we focus on the FEMNIST dataset [77], which is a dataset designed for non-iid federated learning. Section 5.2 presents a generalisation of this setting. More details about our simulations are discussed in Appendix C.

5.1 Non-IID MNIST dataset

Realistic datasets for fully decentralised distributed learning are typically proprietary and not available to the public. A modular benchmarking framework for federated learning is provided by the LEAF benchmark [77], which organises well-established datasets for realistic distributed setting applications. An example is the FEMNIST dataset, which partitions the extended MNIST dataset [84, 85] by the writer of the digit or character into 3550 non-iid sets. The individual sets have, on average, around 227 samples (with a standard deviation of 89 samples) [77].

Due to the small number of samples, we will only train simple models for the data of a single user. Thus, we will run two types of simulations on the FEMNIST dataset. First, we build 101 different models, each being trained on the data of a single different user. Second, we will group users together and train more intricate models on the grouped datasets.

5.1.1 Models for individual users

To begin with, we pick $n = 101$ sub-datasets of the FEMNIST dataset. Each such dataset contains data for 62 distinct classes: 10 digits, 26 lower case letters and 26 upper case letters. We create a binary classification task by labelling all lower case letters by 0, and the upper case letters by 1, while discarding the digits.

The next step is to train 101 different models on these datasets. We deploy three different types of models: random forests (RF) [54], extreme gradient boosting (XGBoost) [86] and light gradient-boosting machine (LGBM) [87]. Each model is trained on 80% of their own data, with the performance on the remaining 20% of the data recorded on the left column of Figure 4. As before, we use accuracy as the performance metric for evaluating the models. Importantly, all models have accuracies better than 50% on their respective validation sets.

The performance of these models can be, of course, slightly improved, by tuning the model parameters using a grid search, for example. Random forests tend to be rather robust to overfitting since the decision trees assembling the forests are independent. As such, the default parameters do typically lead to reasonably well-trained ensembles. While gradient boosters can improve the accuracy of random forests, they can lead to overfitting since they

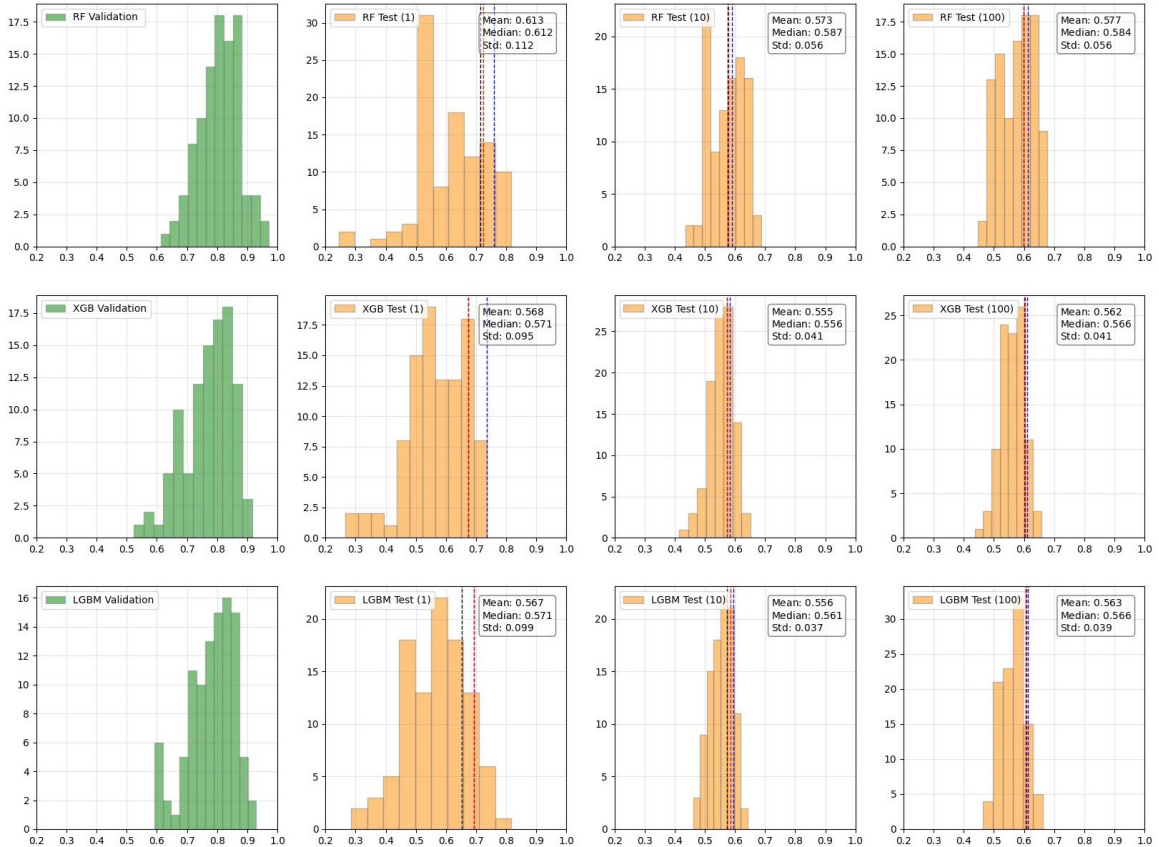


Figure 4: Simulation of Slush consensus learning on FEMNIST dataset for $n = 101$ base learners. *Green*: Distribution of accuracies on validation sets. *Orange*: Distribution of accuracies on test sets using data from 1, 10 or 100 new users, respectively. The dashed lines correspond to: majority ensemble (black), Slush with $\alpha = 6$ (red) and Slush with local α parameters, *i.e.* strong confidence (blue).

repeatedly fit new models to the residuals of previous models. These features can be also seen in the left columns of Figure 4. Nevertheless, these simple models should suffice for our purposes. Note that some of the low accuracies on the test sets can be largely attributed to the small size of the training sets.

For testing, we use three different datasets, from 1, 10 and 100 combined new users. These sets consist of 49, 453 and 4741 samples. The distributions of accuracies in these three testing scenarios are shown in orange in Figure 4, for the three types of models built. It should be noted straightaway that many of the models do not generalise very well to the unseen handwriting. For each of the test datasets, we also generate the majority ensemble, as well as consensus learning ensembles. This will no longer be the case when we increase the size of the training sets in the next subsection.

For the communication phase of the consensus learning algorithms, we run a Slush protocol with a global threshold parameter $\alpha = 6$, and one with local parameters determined by the accuracy on the test sets, in accordance with the concept of strong confidence introduced in Definition 5.1.¹⁰ In both cases, the protocols use $k = 10$ for the sampling of participants. For speeding up the computation time, each communication phase only lasts for 50 rounds per

¹⁰In practical terms, only part of the test subset should be used for determining these local parameters. However, since some of the test sets are rather small, we make an implicit assumption that the test set accuracy would be the same as the accuracy on a smaller subset of this set. This only assumes that this subset is identically distributed to the whole test set.

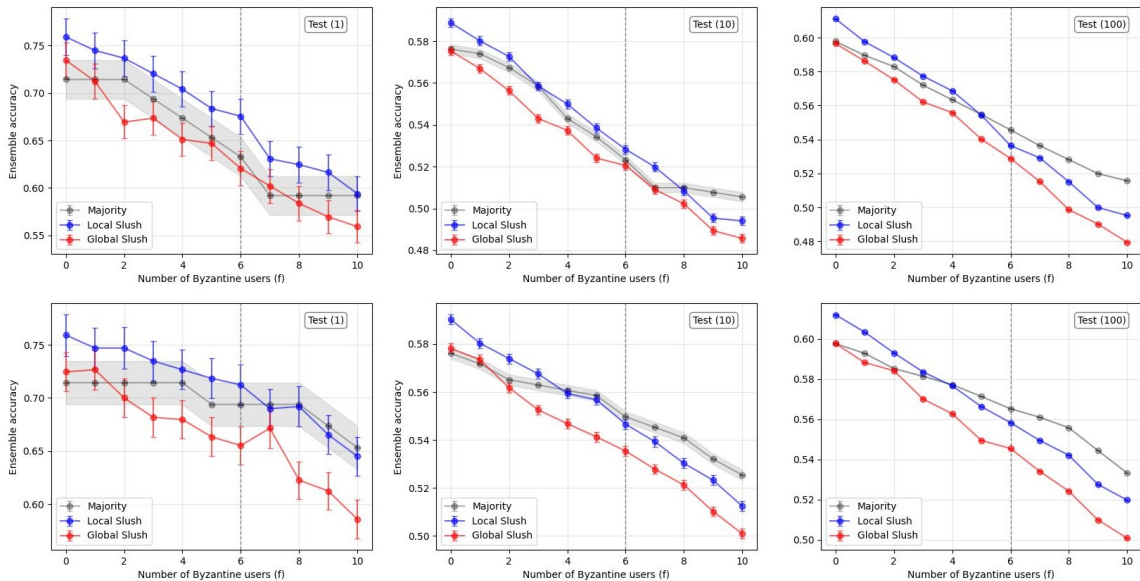


Figure 5: Accuracies of ensembles built from 101 Random Forest models against number of perfectly malicious models. The honest models are trained on non-iid samples from the FEMNIST dataset and tested on data from 1, 10 and 100 new users, respectively. The dashed line $f = 6$ is the value of the α parameter used for the *global* Slush algorithm. *Top row*: strong classifiers turn Byzantine. *Bottom row*: weak classifiers turn Byzantine.

node (so a total of 5000 rounds), which should be more than enough to ensure convergence [27, 67]. The results are indicated with dashed lines in Figure 4.

In Figure 5 we also consider Byzantine participants. For this, f base learners are selected at random from the 101 models and are turned into perfect Byzantine models, as per Definition 4.7. We consider two different samples of base learners that are turned into perfectly Byzantine models: strong classifiers and weak classifiers. More details about the performance of these users on the three test sets can be found in Appendix C.

Figure 5 shows how the accuracies of the ensembles vary when the number of Byzantine participants increases. Let us note that in the context of Byzantine fault tolerance, a centralised majority rule would inevitably have increased Byzantine resilience compared to any consensus protocol. The maximum number of Byzantine participants in the former case is $1/2$ of the total network participants, as opposed to the classical $1/3$ value for consensus protocols [88]. This argument can be extended to consensus learning algorithms. Regardless, consensus learning appears to outperform the majority rule even in the presence of Byzantine users, as long as their number is not too large. The simulations also indicate that when strong models are converted to Byzantine models, the local Slush algorithm has increased resilience as compared to the majority rule.

5.1.2 Models for groups of users

Convolutional neural networks (CNNs) are known to be particularly suitable for image recognition. In this subsection, we train such models on enlarged datasets and subsequently create ensembles of CNNs. To compare the results with the results from the previous section, we will stick to $n = 101$ distinct models. We will use the same test sets as in the previous simulations, which are separate from the training and validation sets.

To avoid overfitting, we use rather simple CNNs, with four hidden layers: a convolutional layer with 8 filters of size (3, 3); a max pooling layer of size (2, 2); a flattening layer; a dense

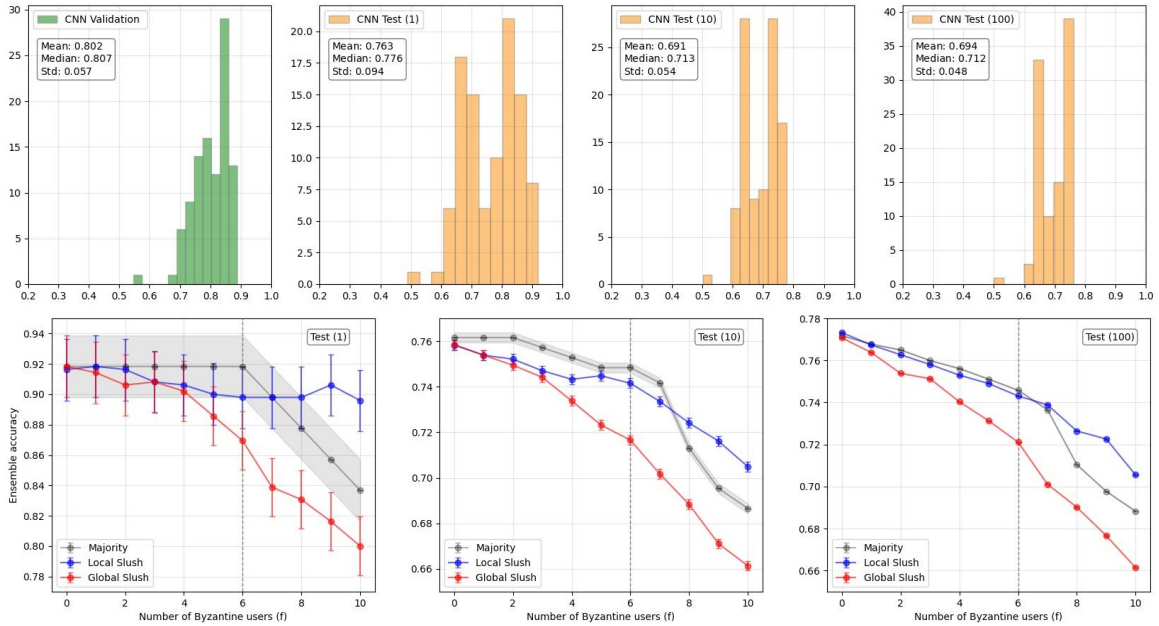


Figure 6: *Top*: Distribution of accuracies on validation sets (green) and test sets (orange) for the 101 CNNs. *Bottom*: Accuracies of ensembles of 101 CNNs against number of perfectly malicious nodes for the three test sets.

layer of 100 neurons with ReLU activation functions and He weight initialisation scheme [89]. For the stochastic gradient descent optimiser, we set the learning rate to 0.01 and use a momentum of 0.9. As the training sets are not too large, we do not use any cross-validation. The models are trained for 10 epochs, with batch sizes of 32 examples.

Perhaps not surprisingly, these CNNs generalise much better than the models trained on data coming from a single user. This is reflected in the distribution of accuracies in the top row of Figure 6. The bottom row in the figure also shows how the ensembles built from these 101 CNNs behave once base learners turn (perfectly) Byzantine. In this case, the sample of 10 base learners that are turned into Byzantine models is chosen to be representative of the whole ensemble of 101 models. More details about this sample are again left to Appendix C. Rather remarkably, we see that the local Slush algorithm still has increased robustness against Byzantine participants and can outperform a centralised majority rule even for larger numbers of Byzantine participants.

5.2 Beta-distributed base learners

The previous simulations on the FEMNIST dataset strengthen our position that consensus learning algorithms can perform better than centralised majority rules if the base learners are diverse. This situation is highly probable within a realistic context and could occur whenever the base learners are trained on non-iid data. These insights are in perfect agreement with our analysis of the control ratio in (4.16).

In this section, we present a generalisation of these results, by generating samples of accuracies for independent base learners from a beta distribution. The beta distribution, $B(a, b)$, is a bounded continuous probability distribution characterised by two real shape parameters, a, b , with $a, b > 0$.¹¹ The distribution approximates a uniform distribution when $a = b = 1$, and a normal distribution for large a, b , with $a \approx b$.

¹¹Here, b should not be confused with the number of blue nodes in the Slush protocol.

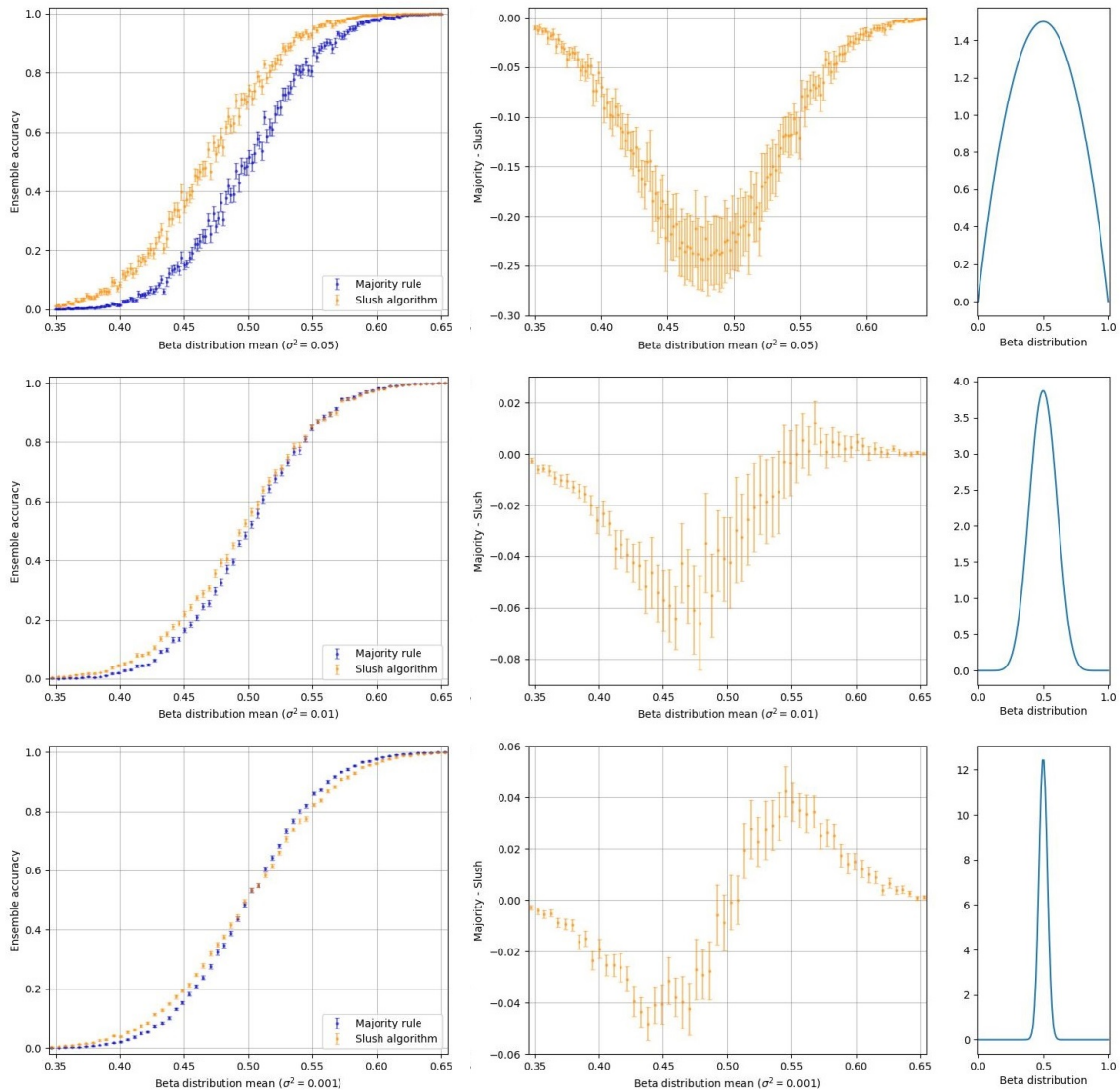


Figure 7: Simulation of Slush consensus algorithm for $n = 101$, $k = 10$, $\mathcal{N}_1 = 100$, $\mathcal{N}_2 = 50$ and local α parameters. *Left column:* Ensemble accuracy for varying mean and fixed variance of the beta distribution. *Middle column:* Difference in accuracy between majority rule and Slush algorithm, for fixed variance. *Right column:* Probability density function of beta distribution with mean equal to 0.5 and variance representative of the row.

It is, of course, difficult to estimate the distribution of accuracies of the base learners without more specifications about the difficulty of the task, or details about the training data used by the base learners. Independent and identically distributed sampling for the training datasets is more likely to lead to distributions that are close to a normal distribution – see *e.g.* [90]. On the other hand, if there are reasonable expectations that the base learners can achieve high accuracies, then a folded normal distribution would be more appropriate.¹² Regardless, our previous simulations do appear to suggest that the beta distribution is a valid candidate.

The simulation can be described by the following steps:

1. A sample of base learner accuracies of size n is generated from a beta distribution with fixed mean and variance.

¹²However, in such scenarios, ensemble methods might not be necessary at all for improving accuracy.

2. For each sample, we generate a vector of individual outputs (*i.e.* a voting profile) and simulate the output of the Slush protocol as well as that of a majority rule using this voting profile. The process is repeated \mathcal{N}_1 for a given sample of accuracies, and the accuracy of the ensemble is determined as the percentage of correctly identified outputs.
3. Step 2 repeats \mathcal{N}_2 times for a given beta distribution, to eliminate any statistical outliers.

The error in the ensemble accuracy is measured as the sampling error for the set of \mathcal{N}_2 values obtained in step 3 of the simulation. The exact value of the error can be determined as discussed in Appendix C, but this would significantly increase the computation time. Finally, the process is repeated for different means and variances of the beta distribution, which is always taken to have a concave density function.

The results of the simulations are shown in Figure 7. As before, we use a Slush algorithm with local threshold parameters α . As expected from the previous theoretical analysis, the Slush algorithm performs significantly better compared to a centralised majority rule for larger variances in the distribution of accuracies.

6 Conclusions and outlook

In this work, we introduced a novel distributed ML paradigm – consensus learning – readily described as a fully decentralised ensemble method that deploys a consensus protocol. We analysed how a typical consensus learning algorithm built on a probabilistic consensus protocol behaves in different scenarios and what improvements it brings to established ML methods. Consensus learning has clear advantages over other distributed learning algorithms, which include communication overhead, resilience against malicious users and protection of private data. Moreover, consensus learning preserves model explainability of typical ensemble weighting methods, offering a high degree of transparency and interpretability.

Our concrete results offer lower bounds on the accuracy of consensus learning classifiers using the Slush consensus protocol, while also describing the behaviour for a large enough number of base learners. We stress again that the proofs of our main results only use *weak* features of the Slush consensus protocol, and could thus be generalised to other probabilistic protocols. In addition to these results, our numerical analysis shows that a relatively simple modification of the Slush consensus protocol can lead to increased Byzantine resilience and a boost in performance.

Our analysis also indicates that a greedy consensus which favours high-accuracy nodes might perform better than Slush algorithms for an ML task. For this purpose, it would be interesting to study federated byzantine agreement (FBA) consensus protocols [91] where each node chooses which participants to trust, as well as hierarchical consensus protocols, which delegate leaders for each communication round. Note, however, that such protocols do typically require precise knowledge of the network, through the so-called *quorum* membership, in order to satisfy safety and liveness guarantees.

Performance boosts can be also achieved by using different local aggregation rules in the communication phase. For a classification task, for instance, these could be weighted majority rules, where, ideally, the weights would be determined by the accuracy of the learners. Such information is typically either not available or, not trustworthy in a distributed setting, but could be stored as some form of past performance in a blockchain implementation. As such, this information would become available to all base learners. Additionally, in such a blockchain implementation, reward mechanisms can create further incentives for participants to be honest, and thus increase the performance of the ensemble. Other Byzantine-resilient aggregation methods (such as those discussed in *e.g.* [9]) could lead to increased robustness

against malicious attacks.

Lastly, consensus learning algorithms can also be applied to other types of ML problems. For regression problems, robust local aggregation rules need to be deployed, similar to Byzantine ML algorithms. The algorithms can also be used for unsupervised ML, similar to other unsupervised ensemble learning methods. We leave a more detailed exposition of these methods for future work.

Acknowledgements

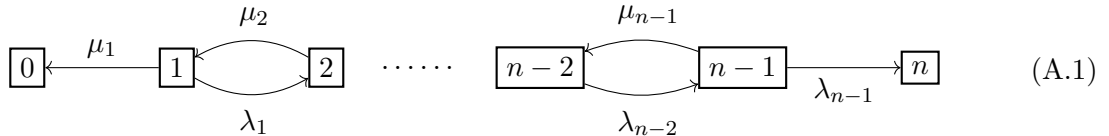
We are very grateful to Charles Grover, Sabrina Guettes and Jernej Tonejc for insightful discussions, constructive feedback and valuable input in many of the proofs presented in this work. We would also like to thank Anne Marie Jackson for helpful comments on the final draft.

A Snow protocols

The Snow family of consensus protocols builds upon the Slush protocol, whose technical aspects we describe in this appendix.

A.1 Slush with honest participants

The Slush protocol is fully described by two control parameters $k \in (0, n]$ and $\alpha \in (k/2, k]$; the former is the size of the sample selected by a node for sending a query, while the latter is the threshold parameter for accepting a query. The protocol can be modelled as a continuous-time Markov chain, with the state \mathcal{S} of the system corresponding to the number of blue nodes (*i.e.* nodes that correctly labelled the transactions) at a given time. This is depicted diagrammatically below:



The process has two absorbing states, all-blue and all-red, respectively. The probability that a query for the blue colour reaches the threshold of α or more votes given b blue nodes in the network can be found from a simple combinatorial exercise, being determined by a hypergeometric distribution,

$$H(n, b, k, \alpha) \equiv \sum_{j=\alpha}^k \binom{b}{j} \binom{n-b}{k-j} \binom{n}{k}^{-1}. \quad (\text{A.2})$$

An important identity for the normalisation of the hypergeometric distribution is

$$\sum_{j=0}^k \binom{n_1}{j} \binom{n_2}{k-j} = \binom{n_1 + n_2}{k}, \quad (\text{A.3})$$

for $n_1, n_2 \in \mathbb{N}$, with $n_1, n_2 \geq k$. Consider the case where all nodes in the network are honest. Then, the Markov chain is described by the following transition rates.

Definition A.1 (Transition rates). *The death μ_b and birth λ_b rates for the state with b blue nodes of the Slush protocol with n nodes and protocol parameters k and α are defined as*

$$\mu_b = bH(n, n-b, k, \alpha) , \quad \lambda_b = (n-b)H(n, b, k, \alpha) . \quad (\text{A.4})$$

The death rate is given by the probability that a given query reaches red consensus; thus, we want one of the b blue nodes to change colour to red. Similarly, for the birth rate, we need one of the $(n-b)$ red nodes to change their colour to blue. These rates satisfy the obvious property $\lambda_b = \mu_{n-b}$. Let us also point out that

$$\lambda_{b < \alpha} = 0 = \mu_{b > n-\alpha} , \quad (\text{A.5})$$

as the voting threshold cannot be reached by the minorities in these cases.

Theorem 6 (Slush absorption, [27]). *Let the configuration of the system at time t be \mathcal{S} , with b blue nodes, where $0 \leq b \leq n$, and $n-b$ red nodes. Then, the probability of absorption \mathcal{R}_b in the all-red state is given by*

$$\mathcal{R}_b = \frac{\sum_{l=1}^{n-b} \prod_{i=1}^{n-l} \mu_i \prod_{j=n-l+1}^{n-1} \lambda_j}{\sum_{l=1}^n \prod_{i=1}^{n-l} \mu_i \prod_{j=n-l+1}^{n-1} \lambda_j} . \quad (\text{A.6})$$

Proof. This is a standard death-birth Markov process, which makes use of the steady-state Kolmogorov equations:

$$(\mu_b + \lambda_b) \mathcal{R}_b = \lambda_b \mathcal{R}_{b+1} + \mu_b \mathcal{R}_{b-1} . \quad (\text{A.7})$$

This recursion relation can be solved explicitly for appropriate boundary conditions, *i.e.* $\mathcal{R}_{b=0} = 1$ and $\mathcal{R}_{b=n} = 0$. See *e.g.* [92], Chapter IV of [93], or Theorem 2 of [27] for explicit proofs. \square

Corollary A.2. *The absorption probability \mathcal{B}_b in the all-blue state given b blue nodes reads*

$$\mathcal{B}_b = \frac{\sum_{l=1}^b \prod_{i=1}^{l-1} \mu_i \prod_{j=l}^{n-1} \lambda_j}{\sum_{l=1}^n \prod_{i=1}^{n-l} \mu_i \prod_{j=n-l+1}^{n-1} \lambda_j} . \quad (\text{A.8})$$

Proof. Since the Markov process has only two absorbing states, we have $\mathcal{B}_b + \mathcal{R}_b = 1$, with \mathcal{R}_b as given in Theorem 6. Using the identity $\lambda_b = \mu_{n-b}$, the result follows. \square

From (A.5), we also note that $\mathcal{B}_b = 0$ for $b < \alpha$. Moreover, as $\lambda_{b < \alpha} = 0$, the sum in the numerator of \mathcal{B}_b in (A.8) will only contribute with terms starting from $l = \alpha$, to $l = b$. Similarly, we notice that $\mathcal{B}_b = 1$ when $b > n - \alpha$, as the absorption probability in the all-red state vanishes in these cases. In practice, however, it might be slightly difficult to work with these expressions.

Proof of Lemma 2.10. An intuitive argument for the first result follows from the neutrality of the protocol, as defined earlier in the context of decisive decision rules: Slush does not discriminate against one of the absorbing states on labelling grounds. The result can also be proved explicitly from the exact expressions of the absorption probabilities, using $\lambda_b = \mu_{n-b}$.

Since $B_{b > \frac{n}{2}} > \frac{1}{2}$, and $B_{b > n-\alpha} = 1$, the second part of the theorem is equivalent to the statement that \mathcal{B}_b is discrete concave for $n - \alpha > b \geq \lceil \frac{n}{2} \rceil$. Thus, we need to show that

$$\mathcal{B}_{b-1} + \mathcal{B}_{b+1} < 2\mathcal{B}_b , \quad (\text{A.9})$$

which, upon using Corollary A.2, amounts to showing that

$$\mu_b < \lambda_b, \quad (\text{A.10})$$

for the above range for b . This can be shown using the explicit form (A.4) of the transition rates. A sufficient condition is given as follows:

$$b \binom{n-b}{j} \binom{b}{k-j} < (n-b) \binom{b}{j} \binom{n-b}{k-j}, \quad (\text{A.11})$$

with $\alpha \leq j \leq k$. Using $b = m + 1 + \delta$, with $0 \leq \delta \leq m$, for $n = 2m + 1$, the above reduces to

$$\prod_{t=k-j}^{j-1} \left(1 - \frac{2\delta + 1}{m + 1 + \delta - t} \right) < 1 - \frac{2\delta + 1}{m + 1 + \delta}, \quad (\text{A.12})$$

which is clearly true.

Finally, when $\alpha = k = 1$, it follows that $\lambda_b = \mu_b$ for any $1 \leq b \leq n - 1$. Then, it is rather straightforward to see that all terms in the summand of (A.8) are equal, and thus $B_b = \frac{b}{n}$. This concludes the proof. \square

Another relevant result used in the proof of Lemma 4.5 is as follows.

Lemma A.3. *Let $F(b^*, n; p)$ be defined as in Lemma 4.5, i.e.*

$$F(b^*, n; p) = \sum_{b=b^*}^n \binom{n}{b} p^b (1-p)^{n-b} = 1 - \sum_{b=0}^{b^*-1} \binom{n}{b} p^b (1-p)^{n-b}, \quad (\text{A.13})$$

for $0 \leq b^* \leq n$. Then,

$$\frac{d}{dp} F(b^*, n; p) = \binom{n}{b^*} b^* p^{b^*-1} (1-p)^{n-b^*}. \quad (\text{A.14})$$

Proof. We prove this statement by induction, as follows. First note that $\frac{d}{dp} F(0, n; p) = 0$ and $\frac{d}{dp} F(1, n; p) = n(1-p)^{n-1}$, in agreement with the above formula. Then, assuming (A.14) holds for $\frac{d}{dp} F(j, n; p)$, we can find $\frac{d}{dp} F(j+1, n; p)$ as follows:

$$\begin{aligned} \frac{d}{dp} F(j+1, n; p) &= - \sum_{b=0}^j \binom{n}{b} p^{b-1} (1-p)^{n-b-1} (b-np) \\ &= \frac{d}{dp} F(j, n; p) - \binom{n}{j} p^{j-1} (1-p)^{n-j-1} (j-np) \\ &= \binom{n}{j} p^{j-1} (1-p)^{n-j-1} (j(1-p) - (j-np)) \\ &= \binom{n}{j} (n-j) p^j (1-p)^{n-j-1} = \binom{n}{j+1} (j+1) p^j (1-p)^{n-j-1}, \end{aligned} \quad (\text{A.15})$$

which concludes the proof by induction. \square

Lemma A.4. *For the Slush protocol with no malicious nodes, the absorption probabilities in the all-blue states satisfy:*

$$\mathcal{B}_{b-1} \mathcal{B}_{b+1} \leq \mathcal{B}_b^2, \quad (\text{A.16})$$

for any b in the range $b \geq \lceil \frac{n}{2} \rceil$. Furthermore, equality occurs only for $b > n - \alpha$.

Proof. Since $b \geq \lceil \frac{n}{2} \rceil$, from the proof of Lemma 2.10 we have $\frac{1}{2}(\mathcal{B}_{b-1} + \mathcal{B}_{b+1}) < \mathcal{B}_b$. By the AM-GM inequality,¹³ we also have:

$$\frac{1}{2}(\mathcal{B}_{b-1} + \mathcal{B}_{b+1}) \geq \sqrt{\mathcal{B}_{b-1}\mathcal{B}_{b+1}}. \quad (\text{A.17})$$

Combining these two yields the required identity. \square

Numerical simulations indicate that this lemma extends to any $1 \leq b \leq n - 1$, with equality for $b < \alpha$ and $b > n - \alpha$. However, in the $b < \lceil \frac{n}{2} \rceil$ range, we now have $\mu_b > \lambda_b$, and thus $\frac{1}{2}(\mathcal{B}_{b-1} + \mathcal{B}_{b+1}) > \mathcal{B}_b$. As such, the argument presented in the above proof no longer applies. Thus, the proof appears to be more intricate for this range, and we leave it for future work. If proved, this extended version of Lemma A.4 would significantly simplify the proof of Theorem 1.

A.2 Byzantine participants

Consider now the scenario described by Theorem 5, *i.e.* there are f Byzantine nodes which follow an ideal adversarial strategy and always respond to a query with the red colour. Let $c = n - f$ be the number of honest participants. In the presence of malicious nodes, the death and birth ratios of the Slush protocol change as follows [27]:

$$\mu_b = b \text{H}(n, n - b, k, \alpha), \quad \lambda_b = (c - b) \text{H}(n, b, k, \alpha), \quad (\text{A.18})$$

for $1 \leq b \leq c - 1$, as $\{0\}$ and $\{c\}$ are the absorbing states, where we also assume $c > f$. Importantly, we have $\mu_{n-b} = \frac{n-b}{c-b}\lambda_b$, while it is still true that $\lambda_b = 0$ for $b < \alpha$, and $\mu_b = 0$ for $b > n - \alpha$. The absorption probabilities $\mathcal{B}_b, \mathcal{R}_b$ in the all-blue and all-red states, respectively, can be found as before, with the distinction that n is replaced by c in (A.8) and (A.6), since the Markov process has $c + 1$ states instead of $n + 1$.

A new and important result in our analysis is the following lemma.

Lemma A.5. *In the Slush protocol with parameters k, α and n nodes, out of which $0 < f < \alpha$ are (perfectly) Byzantine, the absorption probabilities satisfy:*

$$\frac{\mathcal{R}_b}{\mathcal{B}_{n-b}} > \prod_{j=n-b}^b \frac{n-j}{c-j} = \frac{b!}{(n-b-1)!} \times \frac{(c-b-1)!}{(b-f)!}, \quad (\text{A.19})$$

for $\lceil \frac{n}{2} \rceil \leq b \leq n - \alpha$, where $c = n - f$.

Proof. Using the above identities, the numerators of \mathcal{R}_b and \mathcal{B}_{n-b} can be simplified to

$$\begin{aligned} \mathcal{R}_b &\propto \sum_{l=\alpha-f}^{c-b} \prod_{i=1}^{c-l} \mu_i \prod_{j=c-l+1}^{c-1} \lambda_j, \\ \mathcal{B}_{n-b} &\propto \sum_{l=\alpha}^{n-b} \prod_{i=1}^{l-1} \mu_i \prod_{j=l}^{c-1} \lambda_j, \end{aligned} \quad (\text{A.20})$$

as $\lambda_{l < \alpha} = 0$ and $\mu_{l > n - \alpha} = 0$. Due to these identities, and since $n - \alpha < c$, we can rewrite the numerator of \mathcal{R}_b as

$$\mathcal{R}_b \propto \sum_{l=\alpha}^{n-b} \prod_{i=1}^{n-l} \mu_i \prod_{j=n-l+1}^{c-1} \lambda_j. \quad (\text{A.21})$$

¹³Arithmetic and geometric means inequality.

Further algebraic manipulations of this expression lead to

$$\mathcal{R}_b \propto \sum_{l=\alpha}^{n-b} \prod_{i=1}^{l-1} \mu_i \prod_{j=l}^{c-1} \lambda_j \times \left(\prod_{t=l}^{n-l} \frac{n-t}{c-t} \right), \quad (\text{A.22})$$

as long as $\alpha > f$. The fraction inside the last product $\frac{n-t}{c-t}$ is always greater than 1 for $f > 0$ and reaches its minimum for $l = n - b$. Using this value for all terms in the sum, one obtains the required identity. \square

B Proofs of main results

B.1 Proof of Theorem 3

Proof of Theorem 3. To prove the theorem, we proceed as before by computing $\Delta\mathbb{P} = \mathbb{P}_{\text{Maj}}^{(\delta)} - \mathbb{P}_S$, where $\mathbb{P}_{\text{Maj}}^{(\delta)}$ is the majority rule success probability (2.1), but with the starting point of the sum changed to $\lceil \frac{n}{2} \rceil + \delta$. We note that:

$$\Delta\mathbb{P} \leq \sum_{b=\lceil \frac{n}{2} \rceil + \delta}^n \binom{n}{b} p^b (1-p)^{n-b} \mathcal{R}_b - \sum_{b=\lceil \frac{n}{2} \rceil + \delta}^n \binom{n}{b-2\delta} p^{n+2\delta-b} (1-p)^{b-2\delta} \mathcal{B}_{n+2\delta-b}. \quad (\text{B.1})$$

A term-by-term comparison reveals that a set of sufficient conditions for $\Delta\mathbb{P} < 0$ is

$$\frac{\mathcal{R}_b}{\mathcal{B}_{n+2\delta-b}} \binom{n}{b} \binom{n}{b-2\delta}^{-1} \left(\frac{p}{1-p} \right)^{2b-n-2\delta} < 1, \quad (\text{B.2})$$

for $\lceil \frac{n}{2} \rceil + \delta \leq b < n - \alpha$. Each such condition on its own implies that the participant accuracy p must be below a threshold value $p_{\text{th}}(b; \delta) = 1 - (1 + \tau(b; \delta))^{-1}$, with

$$\tau(b; \delta) = \left(\frac{\mathcal{R}_{b-2\delta}}{\mathcal{R}_b} \times \binom{n}{b-2\delta} / \binom{n}{b} \right)^{\frac{1}{2b-n-2\delta}}, \quad (\text{B.3})$$

where we also use Lemma 2.10. When combining these constraints, we are looking for the value of b that minimises $p_{\text{th}}(b; \delta)$.¹⁴ For our purposes, however, this will not be required. Instead, note that $\tau(b; \delta) > 1$, and thus $p_{\text{th}}(b; \delta) > \frac{1}{2}$ for any b and δ . Thus, the true threshold $p_{\text{th}}(\delta)$ will also be larger than $\frac{1}{2}$, as claimed. \square

B.2 Proof of Theorem 1

Proof of Theorem 1. This theorem consists of two separate statements. First, let us look at $\mathbb{P}_S(n, k, \alpha, p) \geq p$, which is equivalent to showing

$$\mathbb{P}_S(n, k, \alpha, p) \geq \frac{1}{n} \sum_{b=0}^n b \binom{n}{b} p^b (1-p)^{n-b}, \quad (\text{B.4})$$

where the RHS is the first moment of the binomial distribution. Next, we split the sums from both sides in two, the first one containing the terms up to $\lfloor \frac{n}{2} \rfloor$, and the second one having

¹⁴This value is only a lower bound on the threshold value $p_{\text{th}}(\delta)$.

the remaining terms. After a simple change of variables, the above condition becomes:

$$\begin{aligned} & \sum_{b=\lceil \frac{n}{2} \rceil}^n \binom{n}{b} \left(\mathcal{B}_b p^b (1-p)^{n-b} + \mathcal{B}_{n-b} p^{n-b} (1-p)^b \right) \\ & \geq \sum_{b=\lceil \frac{n}{2} \rceil}^n \binom{n}{b} \left(\frac{b}{n} p^b (1-p)^{n-b} + \frac{n-b}{n} p^{n-b} (1-p)^b \right). \end{aligned} \quad (\text{B.5})$$

Finally, comparing these expressions term-by-term, and using Lemma 2.10, a sufficient condition for the result to be true is given as:

$$\left(\mathcal{B}_b - \frac{b}{n} \right) \left(\left(\frac{p}{1-p} \right)^{2b-n} - 1 \right) \geq 0, \quad (\text{B.6})$$

for $b > \frac{n}{2}$. This is true by Lemma 2.10, as long as $p > \frac{1}{2}$.

The second statement of Theorem 1 involves the convergence of the Slush algorithm accuracy, \mathbb{P}_S , to unit for large enough n . This statement builds on Lemma 3, leading to a much stronger result. First, according to the aforementioned lemma, large accuracies for the Slush algorithm can occur for the base learner accuracy p in the interval $(q, p_{\text{th}}(q))$. Nevertheless, using the monotonicity of the Slush algorithm proved in Lemma 4.5, we can eliminate the upper bound of this interval.

However, the only remaining issue is to show that $q < p_{\text{th}}(q)$, for any n , such that Lemma 3 can apply. To enlarge the domain of base learner accuracies, we look at $\delta = 1$ or $q = \frac{1}{2} + \frac{1}{n}$. Based on the proof of Theorem 3 – and more precisely on (B.3) – a sufficient condition for $p_{\text{th}}(q) > q$ is

$$\frac{\mathcal{R}_{b-2}}{\mathcal{R}_b} \frac{(b-1)(b-2)}{(n-b+1)(n-b+2)} > \left(\frac{n+2}{n-2} \right)^{2b-n-2}, \quad (\text{B.7})$$

for all $\lceil \frac{n}{2} \rceil + 1 \leq b \leq n - \alpha$. Clearly, all fractions involved are greater or equal than 1. To prove this inequality, consider first the lowest value $b = \lceil \frac{n}{2} \rceil + 1$, when the exponent on the RHS is simply 1; meanwhile, for the LHS we extend Lemma 2.10, such that:

$$\mathcal{R}_{\lceil \frac{n}{2} \rceil} > 1 - \frac{\lfloor \frac{n}{2} \rfloor}{n}, \quad \mathcal{R}_{\lceil \frac{n}{2} \rceil + 1} < 1 - \frac{\lceil \frac{n}{2} \rceil + 1}{n}. \quad (\text{B.8})$$

It follows that:

$$\frac{\mathcal{R}_{\lfloor \frac{n}{2} \rfloor}}{\mathcal{R}_{\lceil \frac{n}{2} \rceil + 1}} > \frac{n+1}{n-3}, \quad (\text{B.9})$$

which is indeed greater than $\frac{n+2}{n-2}$. We would like to present a proof of (B.7) that can easily generalise to other probabilistic consensus protocols. Thus, we want to avoid using too many details that are specific to the Slush protocol. For the Slush protocol, the ratio of the absorption probabilities on the LHS of (B.7) grows (very fast) with b . A detailed argument for this claim is presented below.

Consider the limiting case $\alpha = 1 = k$, where from Lemma 2.10 we have $\mathcal{R}_b = 1 - \frac{b}{n}$ for any $0 \leq b \leq n$. As a result, the ratio $\frac{\mathcal{R}_{b-2}}{\mathcal{R}_b}$ reads $\frac{n-b+2}{n-b}$, which can be shown to be strictly increasing. Additionally, this ratio is strictly larger than $\frac{n+2}{n-2}$ for $b > \lceil \frac{n}{2} \rceil$. When α and k are varied away from this configuration, the rate of change¹⁵ in the absorption probability

¹⁵We define the rate of change as the difference in consecutive absorption probabilities, *i.e.* $\mathcal{R}_b - \mathcal{R}_{b+1}$.

becomes sharper in the region around $b \approx \frac{n}{2}$ and milder otherwise, as depicted in Figure 1. Thus, as long as this is the case, we have $\mathcal{R}_{b-2} - \mathcal{R}_b > \frac{2}{n}$ and thus

$$\frac{\mathcal{R}_{b-2}}{\mathcal{R}_b} > 1 + \frac{2}{n\mathcal{R}_b}. \quad (\text{B.10})$$

The RHS is then larger than $\frac{n+2}{n-2}$ as long as $\mathcal{R}_b < \frac{1}{2} - \frac{1}{n}$, which is of course true for $b > \lceil \frac{n}{2} \rceil$. Importantly, the above argument will hold when the rate of change in the absorption probabilities is larger than $\frac{1}{n}$. However, we are concerned with large values of n , and thus $\frac{1}{n}$ can be made arbitrarily small. The net effect of this is to extend the region where the above argument holds. Finally, for the tail values, *i.e.* when b is large and the rate of change is smaller than $\frac{1}{n}$, the absorption probability \mathcal{R}_b converges to 0, and thus the ratio $\mathcal{R}_{b-2}/\mathcal{R}_b$ diverges.¹⁶

In light of the above reasoning, we can use the ratio $\mathcal{R}_{\lfloor \frac{n}{2} \rfloor} / \mathcal{R}_{\lfloor \frac{n}{2} \rfloor + 1}$ as a placeholder for any value of b in the interval. Then, (B.7) simplifies upon using (B.9) to:

$$\frac{(b-1)(b-2)}{(n-b+1)(n-b+2)} > \left(\frac{n+2}{n-2}\right)^{2b-n-3}, \quad (\text{B.11})$$

for $b \geq \lfloor \frac{n}{2} \rfloor + 2$. An equivalent way of writing this is

$$\frac{(n+2j+1)(n+2j+3)}{(n-2j-1)(n-2j-3)} > \left(\frac{n+2}{n-2}\right)^{2j+2}, \quad (\text{B.12})$$

for $0 \leq j \leq \frac{1}{2}(n-2\alpha-5)$. For $j=0$, the identity can be proved by expanding the terms and comparing the resulting quartic polynomials. More generally, the proof can be done by induction, using

$$\begin{aligned} \frac{(n+2j+1)(n+2j+3)}{(n-2j-1)(n-2j-3)} &= \frac{(n+2j-1)(n+2j+1)}{(n-2j+1)(n-2j-1)} \times \frac{(n+2j+3)(n-2j+1)}{(n-2j-3)(n+2j-1)} \\ &> \left(\frac{n+2}{n-2}\right)^{2j} \times \frac{(n+2j+3)(n-2j+1)}{(n-2j-3)(n+2j-1)} \\ &> \left(\frac{n+2}{n-2}\right)^{2j} \times \left(\frac{n+2}{n-2}\right)^2, \end{aligned} \quad (\text{B.13})$$

where in the last line one proceeds as for the $j=0$ case by expanding the brackets and computing the quartic polynomials. □

B.3 Proof of Theorem 4

Proof of Theorem 4. The distribution for the initial number of blue nodes follows a *Poisson binomial distribution*, given by:

$$\mathbb{P}_{\text{PBN}}(S_n = b) = \sum_{j=\max(0, b-n_2)}^{\min(b, n_1)} \binom{n_1}{j} p_1^j (1-p_1)^{n_1-j} \times \binom{n_2}{b-j} p_2^{b-j} (1-p_2)^{n_2-b+j}. \quad (\text{B.14})$$

¹⁶We are only interested in $b < n - \alpha$, where the strict inequality $\mathcal{R}_{b-2} > \mathcal{R}_b$ holds. Otherwise, $\mathcal{R}_b = 0$ for $b > n - \alpha$.

When $p_2 = 0$, the second part of the expression should be neglected, and the sum reduces to a single term $j = b$. For now, consider generic values of p_1, p_2 . As in the proof of Theorem 2, we do a term-by-term analysis, leading to the control ratio:

$$\kappa_b = \frac{\mathbb{P}_{\text{PBN}}(S_n = b)}{\mathbb{P}_{\text{PBN}}(S_n = n - b)}, \quad \text{for } m + 1 \leq b \leq n - \alpha, \quad (\text{B.15})$$

where $n = 2m + 1$. We are interested in finding the values of p_1 for which $\kappa_b = 1$ and below which $\kappa_b < 1$, for all b . The largest such value occurs in the limiting case $p_2 = 0$ when we have:

$$p_{1,\text{max}}(b) = \left(1 + \tau_b^{\frac{1}{2b-n}}\right)^{-1}. \quad (\text{B.16})$$

Here we introduced:

$$\tau_b = \binom{n_1}{b} \binom{n_1}{b - n_2}^{-1}, \quad (\text{B.17})$$

with $m + 1 \leq b \leq n_1$. Note also that $0 \leq \tau_b \leq 1$ for any b in the given range, and thus $p_{1,\text{max}} > \frac{1}{2}$. The result of the theorem is obtained using $b = \lceil \frac{n}{2} \rceil$ in the previous expressions, which is the value minimising κ_b . To show that this is indeed the case, let $\lambda(b) = \tau_b^{\frac{1}{2b-n}}$, and note that:

$$\left(\frac{\lambda(b)}{\lambda(b+1)}\right)^{(2b-n)(2b-n+2)} = \left(\frac{b+1}{b+1-n_2} \frac{n-b}{n_1-b}\right)^{2b-n} \binom{n_1}{b}^2 \binom{n_1}{b-n_2}^{-2}, \quad (\text{B.18})$$

which is clearly greater than 1 for any $b \geq \lceil \frac{n}{2} \rceil$ as long as $n_2 > 0$. Thus, $\lambda(b)$ is strictly decreasing, as claimed.

Finally, for $p_2 = \frac{1}{2}$, one can check that $p_1 = \frac{1}{2}$ ensures that $\kappa_b = 1$. This follows from the identity (A.3). Thus, for more general $p_2 \in (0, \frac{1}{2})$, the sought-after value of p_1 will lie between the above two limiting cases. □

C Simulation details

In this appendix, we discuss certain aspects of the simulations not covered in Section 5.

Training sets distributions. The extended FEMNIST dataset consists of handwritten digits, lower case letters and upper case letters, partitioned by the writer of the digit or character. To formulate a binary classification problem, we discard the digits and only consider lower and upper case letters. This, of course, reduces the size of the datasets of individual writers. The distribution of dataset sizes used for training the simple ensemble methods (RF, XGBoost and LGBM models) can be seen on the left diagram of Figure 8.

For training CNNs, we group the data of the individual users to create 101 larger datasets. In the LEAF benchmark [77], the 3550 non-iid sets are split into 36 files of 100 users each, with the 50 users in the last file. For simplicity, we split each such file into three groups of equal numbers of users; more precisely, we have groups of 33, 33 and 34 users. The corresponding datasets are then combined in larger sets. The sizes of these training sets are also shown in Figure 8. Note that the test datasets are separate from these training datasets.

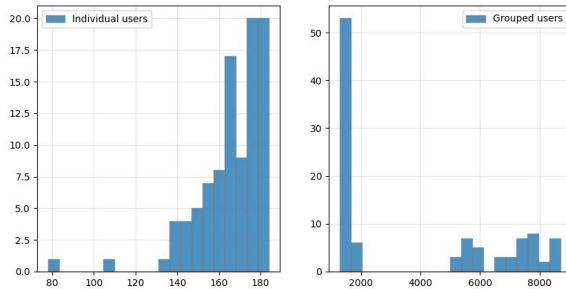


Figure 8: Distribution of dataset size used for training the 101 models. *Left*: datasets of 101 individual users. *Right*: datasets of 101 groups of users deployed for building CNNs.

Selecting Byzantine participants. The simulations on the FEMNIST dataset described in Section 5 do also include Byzantine participants. Consider first the 101 models trained on data from 101 distinct users. For this first scenario, we choose models at random, which are then turned into perfect Byzantine models, *i.e.* their outputs are always incorrect. We then evaluate the accuracy of the majority ensemble and of the Slush algorithms.

The results shown in Figure 5 consider two particular samples of base learners turned into perfect Byzantine models: strong classifiers and weak classifiers. By this we mean that the performances of the classifiers on the test sets are above (below, respectively) the average. These statistics are described in Table 4 for samples of 10 users, and should be compared with

Statistics	Strong classifiers			Weak classifiers		
	Test (1)	Test (10)	Test (100)	Test (1)	Test (10)	Test (100)
Mean	0.699	0.610	0.611	0.601	0.552	0.554
Median	0.714	0.620	0.619	0.531	0.521	0.523
Std	0.091	0.037	0.037	0.114	0.057	0.053

Table 4: Statistics for the two (random) samples of 10 base learners (Random Forests) turned into Byzantine models. Here, each model is trained on data coming from a single user.

these indicated on the top row of Figure 4 for the whole ensemble of 101 classifiers.

A similar procedure is applied to the second simulation, where base learners are trained on data obtained from multiple users. There, we only consider a sample of 10 base learners that is representative of the whole ensemble of 101 models. Their statistics are shown in Table 5, while the ensemble specifications are shown on the top row of Figure 6.

Statistics	Test (1)	Test (10)	Test (100)
Mean accuracy	0.800	0.711	0.715
Median accuracy	0.816	0.726	0.727
Standard deviation	0.090	0.044	0.044

Table 5: Statistics for the sample of 10 base learners (CNNs) turned into Byzantine models.

Error measurement. The plots shown in Section 5 typically include the estimates for the error in ensemble accuracy. The first error to consider is due to the limited number of samples n_{samples} of the test set, being explicitly given by $\varepsilon = (n_{\text{samples}})^{-1}$. This error manifests for any single accuracy estimate of the Slush algorithm. To reduce this error, the communication

phase is repeated $\mathcal{N} = 10$ times for the testing on the datasets coming from 1 and 10 users. Meanwhile, for the test set of 100 grouped users, n_{samples} is large enough for this error to be negligible.

To see how this repetition affects the overall error in the estimate for accuracy of the Slush algorithm we will proceed as follows. Let us model the ensemble accuracy as a random variable Y , whose variance σ is closely related to the previously mentioned error ε through $\sigma \propto \varepsilon^2$. The exact relation is not relevant for our purposes. The next step is to look at the mean of \mathcal{N} such random variables, whose variance becomes:

$$\sigma_{\bar{Y}} = \frac{1}{\mathcal{N}^2} \text{Var} \left(\sum_{i=1}^{\mathcal{N}} Y_i \right). \quad (\text{C.1})$$

If the Y_i variables are uncorrelated, then we have $\sigma_{\bar{Y}} = \sigma/\mathcal{N}$ and thus the error would reduce by a factor of $\sqrt{\mathcal{N}}$. However, these variables are typically not uncorrelated. In the opposite case where the variables are perfectly correlated, we have instead $\sigma_{\bar{Y}} = \sigma$. This latter case applies to the majority rule, which is instead deterministic.

More generally, the error in Slush ensemble accuracy will be smaller than ε , but larger than $\varepsilon/\sqrt{\mathcal{N}}$, and can be found by explicitly computing the covariance matrix of the Y_i variables.

Beta distribution generalities. The mean and variance of a random variable $Z \sim \text{B}(a, b)$ can be expressed as:

$$\mathbb{E}[Z] = \frac{a}{a+b}, \quad \text{var}(Z) = \frac{ab}{(a+b)^2(1+b+1)}, \quad (\text{C.2})$$

which can be easily inverted. As a result, an alternative way of fully specifying the beta distribution is through its mean and variance. Note that $\sigma^2 \equiv \text{var}(Z)$ is limited to the interval $(0, 0.25)$.

References

- [1] J. Konečný, B. McMahan, and D. Ramage, *Federated optimization: Distributed optimization beyond the datacenter*, 2015. arXiv: [1511.03575](#) [[cs.LG](#)].
- [2] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Artificial intelligence and statistics*, PMLR, 2017, pp. 1273–1282.
- [3] M. Abadi *et al.*, *Tensorflow: Large-scale machine learning on heterogeneous distributed systems*, 2016. arXiv: [1603.04467](#) [[cs.DC](#)].
- [4] T. Ben-Nun and T. Hoefler, “Demystifying parallel and distributed deep learning: An in-depth concurrency analysis,” *ACM Computing Surveys (CSUR)*, vol. 52, no. 4, pp. 1–43, 2019.
- [5] J. Verbraeken, M. Wolting, J. Katzy, J. Kloppenburg, T. Verbelen, and J. S. Rellermeier, “A survey on distributed machine learning,” *ACM Comput. Surv.*, vol. 53, no. 2, Mar. 2020, ISSN: 0360-0300. DOI: [10.1145/3377454](#).
- [6] R. Bommasani *et al.*, *On the Opportunities and Risks of Foundation Models*, 2022. arXiv: [2108.07258](#) [[cs.LG](#)].
- [7] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, *Federated Learning: Strategies for Improving Communication Efficiency*, 2017. arXiv: [1610.05492](#) [[cs.LG](#)].

- [8] L. Lamport, R. Shostak, and M. Pease, “The Byzantine Generals Problem,” *ACM Trans. Program. Lang. Syst.*, vol. 4, no. 3, pp. 382–401, Jul. 1982, ISSN: 0164-0925. DOI: [10.1145/357172.357176](https://doi.org/10.1145/357172.357176).
- [9] R. Guerraoui, N. Gupta, and R. Pinot, “Byzantine Machine Learning: A primer,” *ACM Comput. Surv.*, Aug. 2023, ISSN: 0360-0300. DOI: [10.1145/3616537](https://doi.org/10.1145/3616537).
- [10] V. Shejwalkar and A. Houmansadr, “Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning,” in *NDSS*, 2021.
- [11] D. Bouhata, H. Moumen, J. A. Mazari, and A. Bounceur, *Byzantine Fault Tolerance in Distributed Machine Learning : a Survey*, 2022. arXiv: [2205.02572](https://arxiv.org/abs/2205.02572) [cs.DC].
- [12] J. Shi, W. Wan, S. Hu, J. Lu, and L. Y. Zhang, *Challenges and Approaches for Mitigating Byzantine Attacks in Federated Learning*, 2022. arXiv: [2112.14468](https://arxiv.org/abs/2112.14468) [cs.CR].
- [13] A. Cheu, A. Smith, and J. Ullman, *Manipulation attacks in local differential privacy*, 2019. arXiv: [1909.09630](https://arxiv.org/abs/1909.09630) [cs.DS].
- [14] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller, “Inverting gradients - How easy is it to break privacy in 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. 16 937–16 947.
- [15] L. Zhu, Z. Liu, and S. Han, “Deep leakage from gradients,” in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32, Curran Associates, Inc., 2019.
- [16] M. Fredrikson, S. Jha, and T. Ristenpart, “Model inversion attacks that exploit confidence information and basic countermeasures,” in *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, 2015, pp. 1322–1333.
- [17] S. Warnat-Herresthal *et al.*, “Swarm learning for decentralized and confidential clinical machine learning,” *Nature*, vol. 594, no. 7862, pp. 265–270, 2021.
- [18] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, *Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers*. Now Foundations and Trends, 2011, p. 128. DOI: [10.1561/22000000016](https://doi.org/10.1561/22000000016).
- [19] E. M. El-Mhamdi, S. Farhadkhani, R. Guerraoui, A. Guirguis, L.-N. Hoang, and S. Rouault, “Collaborative learning in the jungle (decentralized, byzantine, heterogeneous, asynchronous and nonconvex learning),” in *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, Eds., vol. 34, Curran Associates, Inc., 2021, pp. 25 044–25 057.
- [20] H. Woisetschläger, A. Isenko, S. Wang, R. Mayer, and H.-A. Jacobsen, *A Survey on Efficient Federated Learning Methods for Foundation Model Training*, 2024. arXiv: [2401.04472](https://arxiv.org/abs/2401.04472) [cs.LG].
- [21] O. Sagi and L. Rokach, “Ensemble learning: A survey,” *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 8, no. 4, e1249, 2018.
- [22] S. J. Pan and Q. Yang, “A survey on Transfer Learning,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, 2010. DOI: [10.1109/TKDE.2009.191](https://doi.org/10.1109/TKDE.2009.191).
- [23] S. P. Singh, “Transfer of learning by composing solutions of elemental sequential tasks,” *Machine learning*, vol. 8, pp. 323–339, 1992.
- [24] S. Thrun and L. Pratt, “Learning to learn: Introduction and overview,” in *Learning to learn*, Springer, 1998, pp. 3–17.
- [25] J. Baxter, “A model of Inductive Bias Learning,” *Journal of Artificial Intelligence Research*, vol. 12, pp. 149–198, Mar. 2000, ISSN: 1076-9757. DOI: [10.1613/jair.731](https://doi.org/10.1613/jair.731).
- [26] S. Nakamoto, “Bitcoin: A peer-to-peer electronic cash system,” May 2009.
- [27] T. Rocket, M. Yin, K. Sekniqi, R. van Renesse, and E. G. Sirer, *Scalable and Probabilistic Leaderless BFT Consensus through Metastability*, 2020. arXiv: [1906.08936](https://arxiv.org/abs/1906.08936) [cs.DC].

- [28] T. Hospedales, A. Antoniou, P. Micaelli, and A. Storkey, “Meta-Learning in Neural Networks: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 09, pp. 5149–5169, Sep. 2022, ISSN: 1939-3539. DOI: [10.1109/TPAMI.2021.3079209](https://doi.org/10.1109/TPAMI.2021.3079209).
- [29] L. Metz, N. Maheswaranathan, B. Cheung, and J. Sohl-Dickstein, *Meta-learning update rules for unsupervised representation learning*, 2019. arXiv: [1804.00222](https://arxiv.org/abs/1804.00222) [cs.LG].
- [30] A. Strehl and J. Ghosh, “Cluster ensembles – a knowledge reuse framework for combining multiple partitions,” *Journal of machine learning research*, vol. 3, no. Dec, pp. 583–617, 2002.
- [31] N. Guha, A. Talwalkar, and V. Smith, *One-Shot Federated Learning*, 2019. arXiv: [1902.11175](https://arxiv.org/abs/1902.11175) [cs.LG].
- [32] J. Zhang *et al.*, “Dense: Data-free one-shot federated learning,” in *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., vol. 35, Curran Associates, Inc., 2022, pp. 21 414–21 428.
- [33] D. Dolev, N. A. Lynch, S. S. Pinter, E. W. Stark, and W. E. Weihl, “Reaching approximate agreement in the presence of faults,” *J. ACM*, vol. 33, no. 3, pp. 499–516, May 1986, ISSN: 0004-5411. DOI: [10.1145/5925.5931](https://doi.org/10.1145/5925.5931).
- [34] G. Hinton, O. Vinyals, and J. Dean, *Distilling the knowledge in a neural network*, 2015. arXiv: [1503.02531](https://arxiv.org/abs/1503.02531) [stat.ML].
- [35] H. Chang, V. Shejwalkar, R. Shokri, and A. Houmansadr, *Cronus: Robust and heterogeneous collaborative learning with black-box knowledge transfer*, 2019. arXiv: [1912.11279](https://arxiv.org/abs/1912.11279) [stat.ML].
- [36] D. Li and J. Wang, *Fedmd: Heterogenous federated learning via model distillation*, 2019. arXiv: [1910.03581](https://arxiv.org/abs/1910.03581) [cs.LG].
- [37] C. Roux, M. Zimmer, and S. Pokutta, *On the byzantine-resilience of distillation-based federated learning*, 2024. arXiv: [2402.12265](https://arxiv.org/abs/2402.12265) [cs.LG].
- [38] B. Liu, Z. Ding, and C. Lv, “Distributed training for multi-layer neural networks by consensus,” *IEEE transactions on neural networks and learning systems*, vol. 31, no. 5, pp. 1771–1778, 2019.
- [39] B. Liu and Z. Ding, “Distributed heuristic adaptive neural networks with variance reduction in switching graphs,” *IEEE Transactions on Cybernetics*, vol. 51, no. 7, pp. 3836–3844, 2019.
- [40] S. M. Kazemi *et al.*, “Representation Learning for dynamic graphs: A survey,” vol. 21, no. 1, Jan. 2020, ISSN: 1532-4435.
- [41] F. Bravo-Marquez, S. Reeves, and M. Ugarte, “Proof-of-Learning: A blockchain consensus mechanism based on Machine Learning competitions,” in *2019 IEEE International Conference on Decentralized Applications and Infrastructures (DAPPCON)*, 2019, pp. 119–124. DOI: [10.1109/DAPPCON.2019.00023](https://doi.org/10.1109/DAPPCON.2019.00023).
- [42] Y. Liu, Y. Lan, B. Li, C. Miao, and Z. Tian, “Proof of Learning (PoLe): Empowering neural network training with consensus building on blockchains,” *Computer Networks*, vol. 201, p. 108 594, 2021, ISSN: 1389-1286. DOI: [10.1016/j.comnet.2021.108594](https://doi.org/10.1016/j.comnet.2021.108594).
- [43] Y. Lu, X. Huang, K. Zhang, S. Maharjan, and Y. Zhang, “Blockchain Empowered Asynchronous Federated Learning for Secure Data Sharing in Internet of Vehicles,” *IEEE Transactions on Vehicular Technology*, vol. 69, no. 4, pp. 4298–4311, 2020. DOI: [10.1109/TVT.2020.2973651](https://doi.org/10.1109/TVT.2020.2973651).
- [44] H. Chai, S. Leng, Y. Chen, and K. Zhang, “A Hierarchical Blockchain-Enabled Federated Learning Algorithm for Knowledge Sharing in Internet of Vehicles,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 7, pp. 3975–3986, 2021. DOI: [10.1109/TITS.2020.3002712](https://doi.org/10.1109/TITS.2020.3002712).
- [45] T. G. Dietterich, “Ensemble methods in machine learning,” in *Multiple Classifier Systems*, Berlin, Heidelberg: Springer Berlin Heidelberg, 2000, pp. 1–15, ISBN: 978-3-540-45014-6.
- [46] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *International conference on machine learning*, PMLR, 2017, pp. 1126–1135.
- [47] A. Fallah, A. Mokhtari, and A. Ozdaglar, “Personalized Federated Learning with theoretical guarantees: A model-agnostic meta-learning approach,” 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. 3557–3568.

- [48] P. Bartlett, Y. Freund, W. S. Lee, and R. E. Schapire, “Boosting the margin: a new explanation for the effectiveness of voting methods,” *The Annals of Statistics*, vol. 26, no. 5, pp. 1651–1686, 1998. DOI: [10.1214/aos/1024691352](https://doi.org/10.1214/aos/1024691352).
- [49] T. G. Dietterich *et al.*, “Ensemble learning,” *The handbook of brain theory and neural networks*, vol. 2, no. 1, pp. 110–125, 2002.
- [50] R. Polikar, “Ensemble based systems in decision making,” *IEEE Circuits and Systems Magazine*, vol. 6, no. 3, pp. 21–45, 2006. DOI: [10.1109/MCAS.2006.1688199](https://doi.org/10.1109/MCAS.2006.1688199).
- [51] H. Blockeel, “Hypothesis Space,” in *Encyclopedia of Machine Learning*, C. Sammut and G. I. Webb, Eds. Boston, MA: Springer US, 2010, pp. 511–513, ISBN: 978-0-387-30164-8. DOI: [10.1007/978-0-387-30164-8_373](https://doi.org/10.1007/978-0-387-30164-8_373).
- [52] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan, “A theory of learning from different domains,” *Machine learning*, vol. 79, pp. 151–175, 2010.
- [53] L. Breiman, “Bagging predictors,” *Machine learning*, vol. 24, pp. 123–140, 1996.
- [54] L. Breiman, “Random forests,” *Machine learning*, vol. 45, pp. 5–32, 2001.
- [55] R. E. Schapire and Y. Singer, “Improved boosting algorithms using confidence-rated predictions,” in *Proceedings of the eleventh annual conference on Computational learning theory*, 1998, pp. 80–91.
- [56] D. H. Wolpert, “Stacked generalization,” *Neural networks*, vol. 5, no. 2, pp. 241–259, 1992.
- [57] L. Hansen and P. Salamon, “Neural network ensembles,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 12, no. 10, pp. 993–1001, 1990. DOI: [10.1109/34.58871](https://doi.org/10.1109/34.58871).
- [58] E. H. Frank, *Regression modeling strategies with applications to linear models, logistic and ordinal regression, and survival analysis*, 2015.
- [59] K. K. Ladha, “The Condorcet Jury Theorem, free speech, and correlated votes,” *American Journal of Political Science*, vol. 36, no. 3, pp. 617–634, 1992, ISSN: 00925853, 15405907.
- [60] W. Hoeffding, “On the distribution of the number of successes in independent trials,” *The Annals of Mathematical Statistics*, vol. 27, no. 3, pp. 713–721, 1956. DOI: [10.1214/aoms/1177728178](https://doi.org/10.1214/aoms/1177728178).
- [61] Y. Freund and R. E. Schapire, “A decision-theoretic generalization of on-line learning and an application to boosting,” in *European conference on computational learning theory*, Springer, 1995, pp. 23–37.
- [62] S. Nitzan and J. Paroush, “Optimal decision rules in uncertain dichotomous choice situations,” *International Economic Review*, vol. 23, no. 2, pp. 289–297, 1982, ISSN: 00206598, 14682354.
- [63] M. Castro and B. Liskov, “Practical Byzantine Fault Tolerance and proactive recovery,” *ACM Trans. Comput. Syst.*, vol. 20, no. 4, pp. 398–461, Nov. 2002, ISSN: 0734-2071. DOI: [10.1145/571637.571640](https://doi.org/10.1145/571637.571640).
- [64] A. Avizienis, J.-C. Laprie, B. Randell, and C. Landwehr, “Basic concepts and taxonomy of dependable and secure computing,” *IEEE Transactions on Dependable and Secure Computing*, vol. 1, no. 1, pp. 11–33, 2004. DOI: [10.1109/TDSC.2004.2](https://doi.org/10.1109/TDSC.2004.2).
- [65] C. Cachin, R. Guerraoui, and L. Rodrigues, *Introduction to Reliable and Secure Distributed Programming*, 2nd. Springer Publishing Company, Incorporated, 2011, ISBN: 3642152597.
- [66] C. Cachin and M. Vukolić, *Blockchain Consensus Protocols in the Wild*, 2017. arXiv: [1707.01873 \[cs.DC\]](https://arxiv.org/abs/1707.01873).
- [67] I. Amores-Sesar, C. Cachin, and P. Schneider, *An analysis of avalanche consensus*, 2024. arXiv: [2401.02811 \[cs.DC\]](https://arxiv.org/abs/2401.02811).
- [68] A. Demers *et al.*, “Epidemic algorithms for replicated database maintenance,” in *Proceedings of the sixth annual ACM Symposium on Principles of distributed computing*, 1987, pp. 1–12.
- [69] S. Coretti, A. Kiayias, C. Moore, and A. Russell, “The generals’ scuttlebutt: Byzantine-resilient gossip protocols,” in *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS ’22, Los Angeles, CA, USA: Association for Computing Machinery, 2022, pp. 595–608, ISBN: 9781450394505. DOI: [10.1145/3548606.3560638](https://doi.org/10.1145/3548606.3560638).

- [70] R. Karp, C. Schindelhauer, S. Shenker, and B. Vocking, “Randomized rumor spreading,” in *Proceedings 41st Annual Symposium on Foundations of Computer Science*, 2000, pp. 565–574. DOI: [10.1109/SFCS.2000.892324](https://doi.org/10.1109/SFCS.2000.892324).
- [71] I. Amores-Sesar, C. Cachin, and E. Tedeschi, *When is spring coming? a security analysis of avalanche consensus*, 2022. arXiv: [2210.03423](https://arxiv.org/abs/2210.03423) [cs.DC].
- [72] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A Simple Framework for Contrastive Learning of Visual Representations,” in *Proceedings of the 37th International Conference on Machine Learning*, III, Hal Daume and Singh, Aarti, Ed., ser. Proceedings of Machine Learning Research, vol. 119, PMLR, Jul. 2020, 1597–1607.
- [73] M. A. Kramer, “Nonlinear principal component analysis using autoassociative neural networks,” *AICHE journal*, vol. 37, no. 2, pp. 233–243, 1991.
- [74] J. Rajendran, A. Irpan, and E. Jang, “Meta-Learning Requires Meta-Augmentation,” 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. 5705–5715.
- [75] P. J. Boland, “Majority systems and the Condorcet Jury Theorem,” *Journal of the Royal Statistical Society Series D: The Statistician*, vol. 38, no. 3, pp. 181–189, Dec. 2018, ISSN: 2515-7884. DOI: [10.2307/2348873](https://doi.org/10.2307/2348873).
- [76] M. Fey, “A note on the Condorcet Jury Theorem with supermajority voting rules,” *Social Choice and Welfare*, pp. 27–32, 2003.
- [77] Sebastian Caldas and Sai Meher Karthik Duddu and Peter Wu and Tian Li and Jakub Konečný and H. Brendan McMahan and Virginia Smith and Ameet Talwalkar, *LEAF: A Benchmark for Federated Settings*, 2019. arXiv: [1812.01097](https://arxiv.org/abs/1812.01097) [cs.LG].
- [78] M. J. Fischer, N. A. Lynch, and M. S. Paterson, “Impossibility of distributed consensus with one faulty process,” *Journal of the ACM (JACM)*, vol. 32, no. 2, pp. 374–382, 1985.
- [79] Maksym Zavershynskiy, Medium.com, *Exploring Liveness of Avalanche*, <https://medium.com/@zaver.max/exploring-liveness-of-avalanche-d22f13b2db00>, Accessed: 2024-01-04.
- [80] V. Chvátal, “The tail of the hypergeometric distribution,” *Discrete Mathematics*, vol. 25, no. 3, pp. 285–287, 1979.
- [81] W. Hoeffding, “Probability inequalities for sums of bounded random variables,” *The collected works of Wassily Hoeffding*, pp. 409–426, 1994.
- [82] S. Nitzan and J. Paroush, “Are qualified majority rules special?” *Public Choice*, vol. 42, no. 3, pp. 257–272, 1984.
- [83] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, *Federated Learning with non-IID data*, 2018. arXiv: [1806.00582](https://arxiv.org/abs/1806.00582) [cs.LG].
- [84] Y. LeCun, *The MNIST database of handwritten digits*, <http://yann.lecun.com/exdb/mnist/>, 1998.
- [85] Gregory Cohen and Saeed Afshar and Jonathan Tapson and André van Schaik, *EMNIST: an extension of MNIST to handwritten letters*, 2017. arXiv: [1702.05373](https://arxiv.org/abs/1702.05373) [cs.CV].
- [86] T. Chen *et al.*, “Xgboost: extreme gradient boosting,” *R package version 0.4-2*, vol. 1, no. 4, pp. 1–4, 2015.
- [87] G. Ke *et al.*, “Lightgbm: A highly efficient gradient boosting decision tree,” *Advances in neural information processing systems*, vol. 30, 2017.
- [88] M. Pease, R. Shostak, and L. Lamport, “Reaching Agreement in the Presence of Faults,” *J. ACM*, vol. 27, no. 2, pp. 228–234, Apr. 1980, ISSN: 0004-5411. DOI: [10.1145/322186.322188](https://doi.org/10.1145/322186.322188).
- [89] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org>.
- [90] C. S. Nwosu, S. Dev, P. Bhardwaj, B. Veeravalli, and D. John, “Predicting stroke from electronic health records,” in *2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, 2019, pp. 5704–5707. DOI: [10.1109/EMBC.2019.8857234](https://doi.org/10.1109/EMBC.2019.8857234).

- [91] D. Mazières, *Stellar Consensus Protocol*. *Stellar*, 2021.
- [92] W. Tan, “On the absorption probabilities and absorption times of finite homogeneous birth-death processes,” *Biometrics*, pp. 745–752, 1976.
- [93] S. Karlin and J. McGregor, “The classification of birth and death processes,” *Transactions of the American Mathematical Society*, vol. 86, no. 2, pp. 366–400, 1957.