# Decentral and Incentivized Federated Learning Frameworks: A Systematic Literature Review

Leon Witt, Mathis Heyer, Kentaroh Toyoda, *Member, IEEE,* Wojciech Samek∗, *Member, IEEE,* Dan Li∗

**Abstract**—The advent of Federated Learning (FL) has sparked a new paradigm of parallel and confidential decentralized Machine Learning (ML) with the potential of utilizing the computational power of a vast number of Internet of Things (IoT), mobile, and edge devices without data leaving the respective device, thus ensuring privacy by design. Yet, simple Federated Learning Frameworks (FLF) naively assume an honest central server and altruistic client participation. In order to scale this new paradigm beyond small groups of already entrusted entities towards mass adoption, FLFs must be (i) truly decentralized, and (ii) incentivized to participants. This systematic literature review is the first to analyze FLFs that holistically apply both, blockchain technology to decentralize the process and reward mechanisms to incentivize participation. 422 publications were retrieved by querying 12 major scientific databases. After a systematic filtering process, 40 articles remained for an in-depth examination following our five research questions. To ensure the correctness of our findings, we verified the examination results with the respective authors. Although having the potential to direct the future of distributed and secure Artificial Intelligence, none of the analyzed FLFs is production-ready. The approaches vary heavily in terms of use cases, system design, solved issues, and thoroughness. We provide a systematic approach to classify and quantify differences between FLFs, expose limitations of current works and derive future directions for research in this novel domain.

**Index Terms**—Federated Learning, Blockchain, Incentive Mechanism, Survey.

◆

## 1 INTRODUCTION

CENTRALIZED platforms in the domains of search engines, mobile applications, social media, chat, music and retail have been dominating the respective industries over the past decades. Business models where digital services are exchanged for user data have developed into high-revenue industries with a few single entities controlling the global market within the respective domains [1]. The resulting concentration of user data in a small number of entities, however, poses problems such as the risk of private data leaks [2] or an increasing power imbalance in favor of market-dominating parties [1], [3], [4] which has caused policymakers to enhance data protection for individuals [5]. The need for confidential AI extends beyond B2C markets, such as when entities within the health sector or Internet of Things (IoT) companies are not allowed to collaborate on a common AI model due to sensitive data.

A promising solution that enables the training of Machine Learning (ML) models with improved data security is Federated Learning (FL). In FL, complex models such as Deep Neural Networks (DNNs) are trained in a parallel and distributed fashion on multiple end devices with the training data remaining local at all times. Federated Averaging (`FedAvg`) [6] is a widely applied algorithm for FL where a central authority aggregates a global model from the locally trained models in an iterative process. In theory, FL not only makes previously withheld sensitive data accessible to the machine learning process but also enables efficient training by taking advantage of the ever-increasing computational power of IoT and mobile devices. However, the majority of FL research focuses on advancing the efficiency of the technology, yet incentives and decentralization are necessary requirements for many real-world FL applications, and the prerequisite for FL to evolve from academic research to real-world products with the potential to disrupt the vigorous data and AI industry [7]. In particular, incentives and decentralization address the two major design problems in FL: (i) the star topology of FL that introduces the risk for a single point of failure as well as for authority abuse and prohibits use-cases where equal power among participants is a mandatory requirement, and (ii) the lack of a practical reward system for contributions of participants that hinders this technology from scaling beyond small groups of already entrusted entities towards mass adoption.

Although many proposals of Incentivized and Decentralized Federated Learning Frameworks (FLFs) exist, we have not yet seen any full-fledged production-level FLF. To enhance the development towards production readiness, we compared state-of-the-art solutions despite their heterogeneity in terms of assumptions, use cases, design choices, special focus, and thoroughness by providing a general and holistic comparison framework.

Specifically, we undertake a Systematic Literature

- L. Witt is with the Department of Computer Science, Tsinghua University, Beijing, China, and with the Department of Artificial Intelligence, Fraunhofer Heinrich Hertz Institute, 10587 Berlin, Germany. E-mail: leonmaximilianwitt@gmail.com
- M. Heyer is with Tsinghua University, Beijing, China, and RWTH Aachen University, Aachen, Germany.
- K. Toyoda is with A*STAR, Singapore, and Keio University, Japan.
- W. Samek is with the Department of Electrical Engineering and Computer Science, Technical University of Berlin, 10587 Berlin, with the Department of Artificial Intelligence, Fraunhofer Heinrich Hertz Institute, 10587 Berlin, Germany, and with BIFOLD − Berlin Institute for the Foundations of Learning and Data, 10587 Berlin, Germany. E-mail: wojciech.samek@hhi.fraunhofer.de
- D. Li is with the Department of Computer Science, Tsinghua University, China.
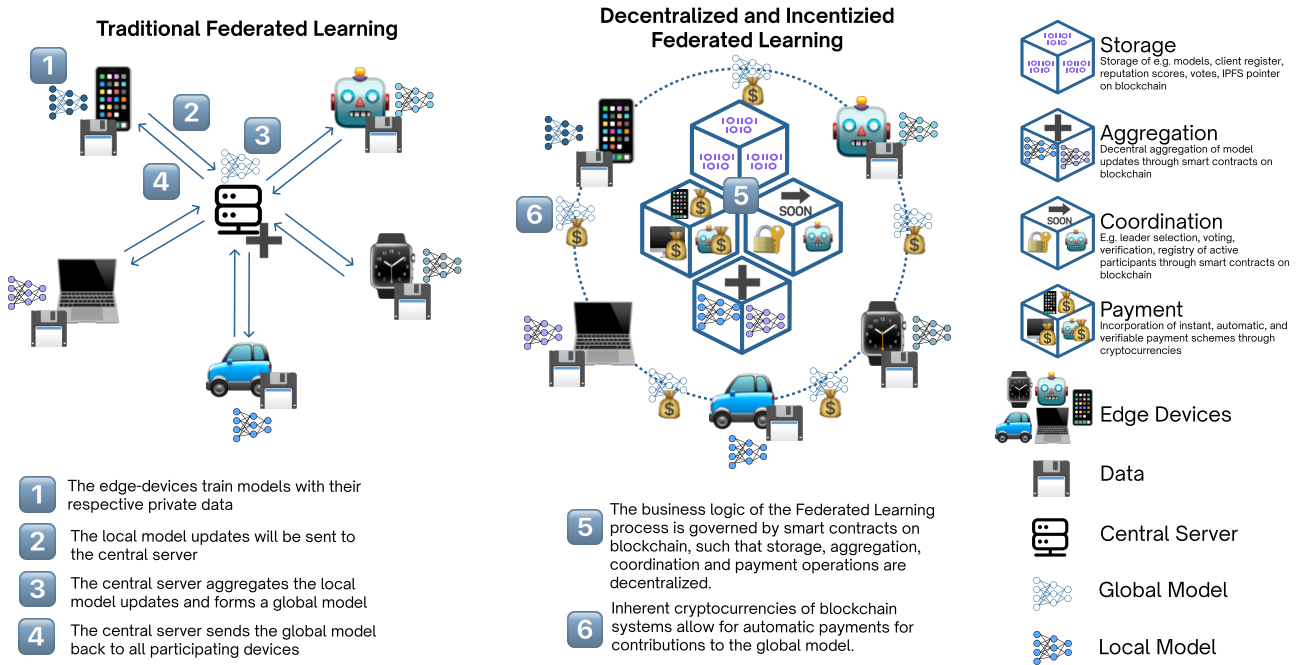
Fig. 1: FL vs. Decentralized and Incentivized FLF.

Review (SLR) examining all relevant articles from twelve scientific databases in the domain of computer science. 422 publications were queried from these databases and filtered for relevant contributions, resulting in 40 papers remaining after three filtering steps. To the best of our knowledge, this is the first comprehensive survey on the design of both decentralized and incentivized federated artificial intelligence systems. The contribution of this paper is threefold:

1) The first comprehensive systematic survey study on the combined topic of decentralized and incentivized FL based on the standardized Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) process, ensuring transparency and reproducibility of the work.
2) A novel comparison framework for an in-depth analysis of incentivized and decentralized FL frameworks which goes beyond existing survey papers by (*i*) pointing out the limitations and assumptions of the chosen game-theoretic approaches, (*ii*) analyzing the existing solutions based on computational and storage overhead on the BC, and (*iii*) an in-depth analysis of the performed experiments.
3) Based on this comparison, we have identified limitations of recent work, derived trade-offs in the design choices, and derived future research directions of decentralized and incentivized FLFs.

The remainder of this paper is structured as follows. In Section 2, we present the technical background of distributed ledger technology and mechanism design in FL systems. In Section 3, we provide an overview of existing surveys on this topic and their respective problem statements. Section 4 summarizes the findings of the

SLR and answers research questions concerning general applications of the FLFs, BC features, Incentive Mechanisms (IMs), and experiments. We derive limitations and further research directions in Section 5. Finally, Section 6 concludes this literature review.

In the Appendix, we outline the details of the methodology, the search process, the selection process, the search terms, and the number of papers retrieved from the respective databases to ensure reproducibility.

## 2 PRELIMINARIES

The following sections discuss the fundamentals of FL, distributed ledger technology, and incentive design, and outline how these technologies are integrated in FLFs.

### 2.1 Federated Learning

FL is a machine learning technique where multiple actors collaboratively train a joint machine learning model locally and in parallel, such that the individual training data does not leave the device as depicted on the left side of Figure 1. This decentralized approach to machine learning was first introduced by Google in 2016 [6] and addresses two key issues of machine learning: (*i*) The high computational effort for model training is relocated from a single central actor to a network of data-owning training devices. (*ii*) As the training data remains on the edge devices, previously inaccessible data of privacy-concerned actors can be integrated into the training process. Thus, "data islands" are prevented.

FL tasks can be classified in horizontal and vertical settings. In horizontal FL, the actors possess the same type of information on different entities, whereas, in vertical FL, the actors have available different information features on the same entity. The latter comes with the

additional challenge of data alignment and information exchange [8]. Furthermore, the FL setting can range from a few collaborating entities, i.e., Cross-Silo (CS), to a federated system of millions of devices, i.e., Cross-Device (CD).

The Federated Averaging (FedAvg) algorithm [6] is a widely adopted optimization algorithm for FL. Its objective is to minimize the empirical risk of the global model $\theta$, such that

$$\arg \min_{\boldsymbol{\theta}} \sum_i \frac{|S_i|}{|S|} f_i(\boldsymbol{\theta}) \tag{1}$$

where for each agent $i$, $f_i$ represents the loss function, $S_i$ is the set of indexes of data points on each client, and $S := \bigcup_i S_i$ is the combined set of indexes of data points of all participants. For that, the calculated gradients from the respective local model training get aggregated in three communication rounds:

1) A central server broadcasts a first model initialization $\theta_{init}$ to a subset of participating clients[1].
2) These clients individually perform iterations of stochastic gradient descent over their local data to improve their respective local models $\theta_i$.
3) In order to create a global model $\theta$, all individual models $\theta_i$ are then sent back to the server, where they are aggregated (e.g., by an averaging operation). This global model is used as the initialization point for the next communication round.

Optimization algorithms for the FL case are open research [7] and variations of FedAvg exist, e.g., FedBoost [9], FedProx [10], FedNova [11], FedSTC [12] and FetchSGD [13].

## 2.2 Blockchain: A Distributed Ledger Technology

A distributed ledger kept by nodes in a peer-to-peer network is referred to as blockchain, first invented by Satoshi Nakamoto through Bitcoin in 2008 [14]. Cryptographic connections of information enable resistance to alteration and immutability. A peer-to-peer consensus mechanism governs the network, obviating the requirement for central coordination [15]. The introduction of general-purpose BCs with smart contract capability supports Turing-completeness [16] and has allowed for the creation of decentralized, immutable, and transparent business logic on top of the BC. Here, smart contracts are computer programs that are decentrally stored on a distributed BC network and automatically executed when a predetermined condition is met.

### 2.2.1 BC data architecture

Even though the data structure varies across different BC, the structure can roughly be categorized into six layers.

- **Data Layer**: The data layer defines how new information will get stored and how the respective blocks are designed. Blocks typically contain a block

1. Clients, workers, and agents are used interchangeably.

header and block body [14], [16]. The block header is a collection of metadata about the block and a summary of the transactions included in the execution payload. In Ethereum, the block body is a bundled unit of information that include an ordered list of transactions and consensus-related information [16].

- **Network Layer**: BCs are peer-to-peer networks with nodes that require defined protocols for communication between them. The protocol comprises exchanging requests and answers between particular nodes (one-to-one communication) as well as "gossiping" information (one-to-many communication) through the network [17]. To make sure they are delivering and receiving the right information, each node must abide by a set of networking regulations.
- **Consensus Layer**: The consensus mechanism refers to the entire stack of protocols, incentives, and ideas that allow a network of nodes to agree on the state of a BC. The two major types of consensus mechanisms for public BCs are Proof of Work (PoW) and Proof of Stake (PoS). PoW is used in Bitcoin [14] and the former version of Ethereum [16], where nodes in the network compete to find a specific hash value below a given number to prove that a certain amount of a specific computational effort has been expended in order to add blocks to the network. The fact that it would take 51% of the network's computer power to commit fraud on the chain ensures the network's security [14]. In PoS-based consensus mechanisms, blocks can be added by randomly chosen validators who have staked significant amounts of cryptocurrencies. The system is designed such that the staked cryptocurrencies get slashed when the validators act maliciously, securing the network crypto-economically [18].
- **Incentive Layer**: To incentivize participation in the consensus mechanism, BC systems contain an inherent cryptocurrency reward for either contributing computational effort (i.e., PoW) or staking cryptocurrencies (i.e., PoS). FL-specific BC systems may reward operations beyond securing the BC such as the storage of ML models [19], [20], [21], the aggregation of gradients [19], [22], [23], [24], [24] or contribution calculations [22], [23], [25].
- **Contract Layer**: Smart contracts are simple programs that run on the respective virtual machine of the BC. The Ethereum Virtual Machine (EVM) is a Turing-complete environment for smart contracts most commonly used across other BCs such as Polygon [26], BNB Chain [27] and Avalanche [28]. A smart contract is a collection of code (its functions) and data (its state) that resides at a specific address on the respective BC.
- **Application Layer**: Decentralized Applications (Dapps) on top of smart contracts cannot be censored, allow for anonymous participation, have zero downtime, and are compatible with other Dapps on the same BC.

### 2.2.2  BC-based FL to ensure equal power

Due to its intrinsic features, distributed ledger technology is capable of mitigating open issues in the FL context, namely:

- **Decentralization**: Workers are subject to a power imbalance and a single point of failure (SPoF) in server-worker topologies. A malicious server might refuse to pay reward payments or exclude employees at will. Furthermore, a server-worker architecture is incompatible with a situation in which numerous entities have a shared and equal stake in the advancement of their respective models. BC technologies' decentralization provides a federal system for entities of equal authority without the need for a central server.

- **Transparency and Immutability**: Data on the BC can only be updated, not erased, as every peer in the system shares the same data. In a FL environment, a clear and immutable reward system ensures worker trust. On the other side, each client is audited, and as a result, can be held accountable for malevolent activity.

- **Cryptocurrency**: Many general-purpose BC systems include cryptocurrency capabilities, such as the ability to incorporate payment schemes within the smart contract's business logic. Workers can be rewarded instantly, automatically, and deterministically based on the FL system's reward mechanism.

Therefore, BC systems [16], [29], [30] have the potential to mitigate the first issue of FL by ensuring trust through their inherent properties of immutability and transparency. They enable decentralized federations to mitigate dependencies on a central authority.

Figure 1 (right) depicts the four major functions BC helps to facilitate in the FL process, namely aggregation, coordination, storage, and payment.

- **Aggregation**: In regular `FedAvg`-based FL, a central authority collects and aggregates the clients' gradients. BC can decentralize the process by performing the aggregation on-chain.

- **Coordination**: Leader selection for the aggregation process, voting, verification, and onboarding of new clients are necessary for real-world FL but undefined steps in contemporary research. BC can provide a trustless and transparent infrastructure for those steps.

- **Storage**: BC provides immutable and transparent storage for information where access is shared among clients.

- **Payment**: BCs inherent functionality of cryptocurrencies allows for automated payment schemes to reward clients for the exerted effort. These four major operations of BC in FL are discussed in detail in Section 4.2.

## 2.3  Incentive Mechanism

FL use cases where pseudo-anonymous clients are expected to participate and invest their data and computational power cannot assume altruism but require compensation in any real-world scenario. Mechanism Design (MD), which is a field of economics, attempts to implement a protocol, system, or rule so that the desired situation (e.g., every participant contributes informed truthful model updates) is realized in a strategic setting, assuming that each participant acts rationally in a game theoretic sense [31].

The purpose of incorporating MD into FL is to incentivize clients to (*i*) put actual effort into obtaining real and high-quality signals (i.e., training the model on local data) and (*ii*) submit model updates truthfully despite not being monitored directly. Such incentives can be distributed using BC infrastructure and their underlying cryptocurrencies. An appropriately designed mechanism ensures a desired equilibrium when every worker acts rationally and in their own best interest. Moreover, a mechanism ideally has low complexity and is self-organizing, avoiding the need for Trusted Execution Environments (TEE) [32] or secure multi-party computation (MPC), yet makes assumptions about the degree of information available.

The process of designing a FL protocol with MD consists of (*i*) designing a mechanism and (*ii*) a theoretical analysis. The former determines the whole procedure of FL including a reward policy. The reward policy defines (*i*) how to measure clients' contribution to the overall model and (*ii*) how to distribute rewards. Measuring contributions in FL is challenging since the aggregated gradients do not reveal explicit information about their effect on the overall performance. Additionally, a myriad of design choices for reward distribution exists, e.g., whether rewards should be given to the top contributor (i.e., winner-takes-all) or to multiple workers where rewards are equally distributed among all contributors or unequally distributed based on the workers' contribution.

### 2.3.1  Theories behind mechanism design

We classify underlying theories for mechanism design broadly into two categories, namely (*i*) game theory and (*ii*) auctions, based on [33].

The theory assumes that the clients' utility $U$ is defined by expected profits $\Pi$ (e.g., prizes) minus costs $C$ (e.g., costs of model training and data collection).

$$U = \Pi - C \qquad (2)$$

Assuming individual rationality, clients choose their actions to maximize their utilities. In this context, the interactions of choices that produce outcomes concerning utilities are referred to as games in the scientific literature. Games can be classified into cooperative and non-cooperative. A non-cooperative game is a game where each client individually determines their strategies so that their utilities are maximized, while a cooperative game maximizes the utility of the group. A game is called imperfect when a client does not know the others' information (e.g., utilities, strategies, etc.). When all the clients know others' information, such a game is called perfect. Popular game-theoretic methods applied in FL are Stackelberg games, contest theory, and contract theory.

**Stackelberg game**: A leader (e.g., a task requester) determines their strategy, and followers (i.e., clients or workers) determine theirs according to the leader's

action [34]. A task requester can be a leader who determines a reward first, and clients determine their effort and how much they should exert based on the condition. Seminal work on Stackelberg games includes, e.g., Khan *et al.* who motivate the modeling of FL as Stackelberg game and propose an IM based on their best response algorithm [35]. Another example of a Stackelberg game-based incentive mechanism in centralized FL settings can be found with Zhang *et al.* [36]. The authors tackle the challenges of information secrecy and contribution measurement by training a deep reinforcement learning-based IM that allows optimal pricing for the central server and optimal training strategies for the participating clients.

**Contest Theory**: Contrary to Stackelberg games, clients need to exert efforts before joining the contest [37]. The process of FL can be seen as a contest as clients must train a model with their own local datasets while they are not guaranteed to receive prizes.

**Contract Theory**: In contract theory, an employer has to agree on a contract with employees given the situation that the employees may claim false capabilities [38]. This could be the case in FL as task requesters do not exactly know clients' capabilities (e.g., cost, computational resources, etc.).

**Auctions**: Auctions are applicable in designing the mechanisms of FLF as they optimally allocate resources such as computational resources or amount of data, based on clients' reports. A task requester posts a FL task, potential clients bid with sealed information such as computational cost and resource, and the requester assigns a FL task to winning clients. There are several auctions to determine winners (e.g., first-price sealed-bid (FPSB) auction, second-price sealed-bid (SPSB) auction, Vickrey-Clarke-Groves (VCG) auction, all-pay auctions) [39]. A detailed analysis of all auction types is beyond the scope of this paper which is limited to the most popular auctions applied in FLFs: a FPSB auction is an auction where no bidders know others' bids and the highest bidder pays the price that they bid. A SPSB auction is similar to the FPSB, but the highest bidder only needs to pay the price that the second highest bid. A VCG auction is a sealed-bid auction for multiple resources. It is designed to achieve socially optimal resource allocation by charging winners of an auction the social loss they cause to others. This prevents clients from bidding their false valuations to win. An all-pay auction is an auction where all bidders need to pay regardless of whether or not they win. A Tullock contest is one of the most famous all-pay auctions [40].

### 2.3.2 Desirable properties

Game theory and auctions provide a strong guarantee that a designed mechanism possesses desirable properties. Zeng *et al.* summarized the main properties that a mechanism should possess in FL, namely incentive compatibility (IC), individual rationality (IR), Pareto efficiency (PE), collusion resistance (CR) and fairness, balanced budget (BB) [8]. IC is fulfilled when entities cannot be better off by deviating from their optimal strategies, and IR refers to the assumption that contributors would not participate if their respective utility was negative as in Eq. (2). A game is PE when it guarantees that the sum of profits is maximized. CR is achieved when no participant

TABLE 1: Comparison of related survey papers.

| Ref. | BC | FL | IM | Experiment Analysis |
|---|---|---|---|---|
| [43] | ✓ | ✓ | | |
| [42] | Partially | ✓ | | |
| [8] | Partially | ✓ | ✓ | |
| [44] | | ✓ | ✓ | |
| [33] | | ✓ | ✓ | |
| [45] | Partially | ✓ | Partially | |
| [46] | ✓ | ✓ | Partially | |
| **This work** | ✓(Detailed) | ✓ | ✓ | ✓ |

can be better off by colluding with others. A game is said to be fair when fairness, e.g., a payoff to the contribution, is preserved [41]. Finally, a game is BB when it is sustainable without external money inflows. Section 4.3 discusses how these theories and properties are adopted in the FLFs' MD.

## 3 RELATED SURVEYS

To the best of our knowledge, this is the first analysis of holistic frameworks for fully decentralized FL with rewards for the participating clients. Yet, we have identified several survey papers in the context of either MD and FL [8], [33], [42] or BC and FL [8], [43], [44]. TABLE 1 shows the comparison of the related survey papers and our own.

Hou *et al.* investigate the state-of-the-art BC-enabled FL methods [43]. They focus on how BC technologies are leveraged for FL and summarized them based on the types of BC (public or private), consensus algorithms, solved issues, and target applications.

The other related survey papers focus on IM for FL [8], [33], [42], [44]. Zhan *et al.* survey the IM design dedicated to FL [42]. They summarize the state-of-the-art research efforts on the measures of clients' contribution, reputation, and resource allocation in FL. Zeng *et al.* also survey the IM design for FL [8]. However, in this publication, the authors focus on IM such as Shapley values, the Stackelberg game, auction, context theory, and reinforcement learning. Besides [42] and [8], Ali *et al.* [44] also summarize involved actors (e.g., number of publishers and workers), evaluation datasets as well as advantages and disadvantages of the mechanisms and security considerations. Tu *et al.* [33] provide a comprehensive review of economic and game theoretic approaches to incentivize data owners to participate in FL. In particular, they cluster applications of Stackelberg games, non-cooperative games, sealed-bid auction models, reverse action models as well as contract and matching theory for incentive MD in FL. Nguyen *et al.* investigate opportunities and challenges of BC-based FL in edge computing [45]. Finally, Wang *et al.* surveyed BC-based FLF with a particular focus on FLF system compositions [46].

As revealed from our analysis and summarized in TABLE 1, the existing survey papers lack a holistic analysis of FL frameworks that are both decentralized *and* incentivized. Such a review is urgently needed since the simultaneous implementation of these features comes with additional interdisciplinary challenges while being crucial to establishing a fair and trustworthy FL framework to the benefit of the data owner. To fill this research gap, this paper
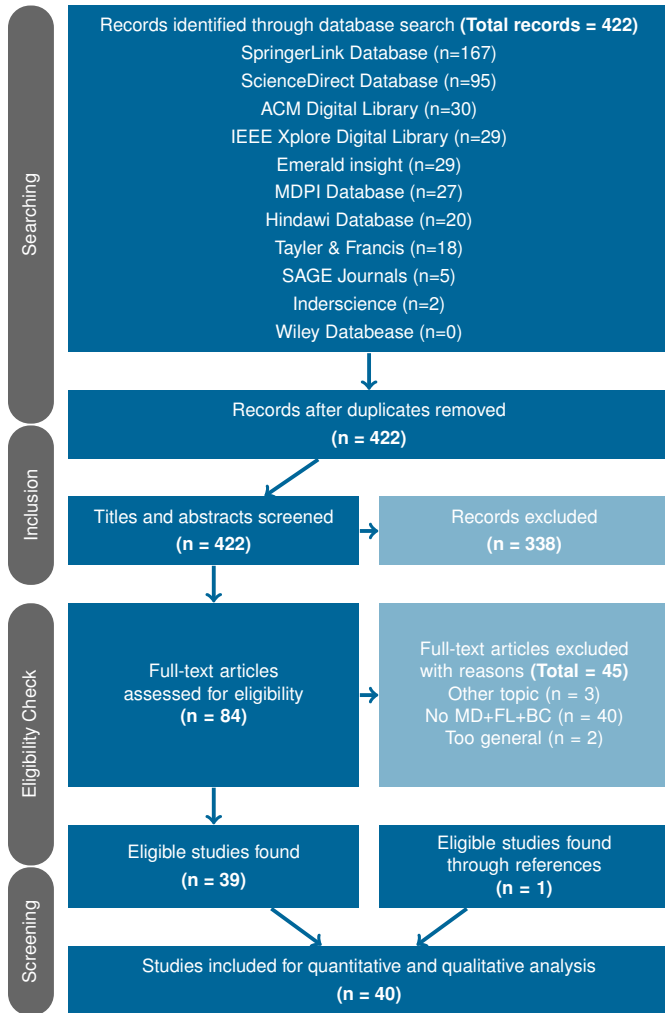
Fig. 2: Flow diagram of the search and screening step in the PRISMA methodology.

provides the first systematic literature review on the topic of BC-enabled decentralized FL with IM.

# 4 SYSTEMATIC LITERATURE REVIEW

The goal of the systematic literature review is the identification of decentralized collaborative learning solutions where participation is rewarded. For that, relevant publications are retrieved, filtered, and analyzed following a methodical, reproducible procedure. The procedure is inspired by the PRISMA methodology [47] and augmented with the guide for information systems proposed by Okoli *et al.* and Kitchenham *et al.* [48], [49]. The five core steps of the systematic approach include (*i*) defining research questions, (*ii*) searching for literature, (*iii*) screening, (*iv*) reviewing, (*v*) selecting and documenting relevant publications, and extracting relevant information.

A detailed visualization of the search and screening step can be found in Figure 2. Details about the SLR are outlined in the Appendix.

Subsections 4.1 to 4.5 systematically present the results of the literature review by answering our five research questions:

RQ1  Overview: (*i*) What are the possible applications of FLF? (*ii*) What problems were solved?

RQ2  BC: (*i*) What is the underlying BC architecture? (*ii*) How is BC applied within the FLF and what operations are performed? (*iii*) Is scalability considered?

RQ3  IM: (*i*) How are IMs analyzed? (*ii*) How are the contributions of workers measured?

RQ4  FL: (*i*) Is the performance of the framework reported? (*ii*) How comprehensive are the experiments? (*iii*) Are non-IID scenarios simulated? (*iv*) Are additional privacy methods applied? (*v*) Is the framework robust against malicious participants?

RQ5  Summary: What are the lessons learned from the review?

Each subsection is complemented by an explanatory table that classifies the considered papers according to categories defined in TABLE 2. Hereafter, we use n.a. (not applicable) for the items that are not applicable. Likewise, we use n.s. (not stated) for the items that should be stated but are not. We also use ✓ for the items that satisfy a condition while leaving cells empty when they do not.

## 4.1 RQ 1: Overview

### 4.1.1 RQ 1-1: What are possible applications of FLF?

Table 3 shows the summary of FL frameworks. Although most of the surveyed papers do not target specific applications (28 out of 40) due to the generalizability of neural networks, some are dedicated to specific applications, namely IoT (6 out of 40), Internet of Vehicles (IoV) (5 out of 40) and Finance (1 out of 40). Applications of FLF in special domains may require additional constraints and characteristics. The heterogeneity of the required properties across those domains leads to vast differences in the design choices of function, operations, and storage of BC, contribution measurement, and privacy requirements.

One of the major application scenarios is IoT (e.g., [51], [52], [63]). Sensor-equipped devices collect environmental information and execute model updates thanks to advances in neural engines while edge servers are often assumed to aggregate models that are trained by local sensor devices. For instance, power consumption measured at smart homes can be used for training an AI model of energy demand forecast [72]. Zhao *et al.* propose a FLF for smart home appliance manufacturers to obtain AI models trained with their customers' usage information [52].

Some solved issues pertaining to the IoT-based FL [51], [68]. Zhang *et al.* propose a FL-based failure device detection method that takes into account the fact that sensor readings are often imbalanced since sensors are, in general, not deployed uniformly in a sensing area [51]. They propose a modified FedAvg algorithm called centroid distance weighted federated averaging (CDW_FedAvg) to obtain accurate models when local datasets are imbalanced at the devices. As sensor devices may not have enough resources to solely train neural networks, it is important to determine whether to delegate computationally intensive tasks to edge servers. Qu *et al.* propose an algorithm to determine whether to offload computation to edge servers when

TABLE 2: Definition of columns in the overview tables.

| Table | Column | Definition | Examples |
|---|---|---|---|
| **FL (TABLE 3)** | Application | Fields of applications | Generic, IoT |
| | Setting | Whether a setting of FL is given | CS, CD |
| | Actors | Actors assumed in the FLF | Workers |
| | Setup | Whether how a system is set up was given (e.g., who deploys a BC) | ✓ |
| | Domains | To which domain each work contributes | |
| | | SPoF: Single point of failure | ✓ |
| | | BC: Blockchain | ✓ |
| | | FL: Federated learning | ✓ |
| | | IM: Incentive mechanisms | ✓ |
| | | CM: Contribution measurement | ✓ |
| | | SP: Security & privacy | ✓ |
| **BC (TABLE 4)** | Operations | Whether the following operations are executed on the BC | |
| | | Agg.: Model aggregation | ✓ |
| | | Cor.: Coordination | ✓ |
| | | Pay.: Payment | ✓ |
| | | Str: Storage | ✓ |
| | BC | BC used in the FLF | Ethereum |
| | Consensus | Consensus mechanism used | PoW |
| | On-chain | Items: Types of data stored on-chain | ✓ |
| | | Eval.: Whether the on-chain storage amount was evaluated | ✓ |
| | Off-chain | Whether off-chain storage (e.g., IPFS) is used | ✓ |
| | Scalability | Whether scalability is considered | ✓ |
| **IM (TABLE 5)** | Sim. | Whether a game was evaluated via simulation | ✓ |
| | Theoretical analysis | Whether a game was theoretically analyzed | ✓ (Contract theory) |
| | Contribution measurement | Costs: Costs considered in analysis | Energy |
| | | Metrics: Metrics to validate workers' contribution | Accuracy |
| | | Abs.: Whether the metric is absolute | Accuracy |
| | | Rel.: Whether the metric is relative | Accuracy |
| | | Rep.: Whether reputation is considered | Accuracy |
| | | Validator: Actors that validate workers' contribution | Task requesters |
| **Experiments (TABLE 6)** | Tasks | Types of ML tasks | |
| | | Clf: Classification | ✓ |
| | | Rgr: Regression | ✓ |
| | Datasets | Datasets used for evaluation | MNIST |
| | #Clients | Number of clients in the experiments | 10 |
| | Algorithms | FL algorithms used | FedAvg |
| | Privacy | Whether privacy protection methods were applied | ✓ |
| | Non-IID | Whether the non-IID condition was assumed | ✓ |
| | Adversaries | Whether security analysis was given against the following attacks | |
| | | BT: Blockchain tampering | ✓ |
| | | RP: Random model poisoning | ✓ |
| | | RT: Reputation tampering | ✓ |
| | | SP: Systematic model poisoning | ✓ |
| | Imp. | Whether the BC part was implemented for evaluation | ✓ |
| | Per. | Whether the performance of FL models was measured | ✓ |

communications between IoT devices and edge computers are unreliable [68]. Beyond, Liu *et al.* [83] reflect on FL in the context of 6G communication and how both technologies are expected to empower each other. The review pinpoints communication cost, security, privacy, and training interference as the key challenges of FL and 6G communication.

FL is beneficial to many scenarios in ITS or IoV, e.g., optimized routing, congestion control, and object detection for autonomous driving ( [20], [59], [64], [74], [76]). Vehicles collect local information and train local models with collected data. Models are often aggregated by devices called Road Side Units (RSUs) and Mobile Edge Computers (MECs) which are often deployed on the road. In IoV, the CD setting is often preferred as mostly the same types of sensors are used to measure road conditions, and thus the common neural network model structure is shared by vehicles. As different locations have different road conditions, users need locally-optimized models, and thus scalability is a key issue. Furthermore, we need extra protection for users' location privacy. Zou *et al.* propose a FLF for a knowledge trading marketplace where vehicles can buy and sell models that vary geographically [74]. Chai *et al.* propose multiple BCs to deal with geographically dependent

TABLE 3: Overview of decentral and incentivized FL frameworks.

| Ref. | Application | Setting | Actors | Setup | Domains SPoF | BC | FL | IM | CM | SP |
|------|-------------|---------|--------|-------|------|----|----|----|----|----|
| [22] | Generic | n.s. | Workers | ✓ | ✓ | | ✓ | ✓ | ✓ | |
| [19] | Generic | CS | Workers | ✓ | ✓ | | | ✓ | | ✓ |
| [50] | Generic | n.s. | Contributors, miners | ✓ | ✓ | | ✓ | | | ✓ |
| [51] | IoT | n.s. | Clients, central organization | | ✓ | | ✓ | ✓ | | ✓ |
| [52] | IoT | n.s. | Clients, miners | ✓ | ✓ | | | | | ✓ |
| [53] | Generic | n.s. | Aggregation servers, workers | ✓ | | | | ✓ | | |
| [54] | Generic | n.s. | Workers, task publishers, miners | | | | | ✓ | | |
| [20] | IoV | n.s. | MEC, MBS, ADV | | | | | ✓ | | ✓ |
| [25] | Generic | n.s. | Model requesters, FL servers, clients | | | ✓ | | ✓ | | |
| [55] | Generic | CD | Administrator, requesters, workers, validators | ✓ | | | | ✓ | | ✓ |
| [21] | Generic | n.s. | Central server, workers | | | | | ✓ | | |
| [56] | Generic | CS | Workers, leaders, aggregation server | | | | | ✓ | | ✓ |
| [57] | Generic | n.s. | Edge devices, fog nodes, cloud | ✓ | ✓ | | | | | |
| [58] | Generic | CS | Users, edge devices, cloud | | ✓ | | | | | ✓ |
| [59] | IoV | CD | Vehicles, RSUs, BSs | | ✓ | | | | | |
| [60] | Generic | CD | Task publishers, workers | | | | | ✓ | | |
| [61] | Generic | n.s. | Trainers, buyers, reporters, data processors | ✓ | ✓ | | | ✓ | | |
| [62] | Generic | CS | Data owners | ✓ | | | | ✓ | | |
| [63] | IoT | CD | Local devices, BSs, MEC nodes | ✓ | ✓ | | | ✓ | | |
| [64] | IoV | CD | Vehicle edge nodes, BC nodes | ✓ | ✓ | | ✓ | ✓ | | |
| [65] | Generic | n.s. | Trainer nodes | ✓ | ✓ | | | ✓ | | ✓ |
| [66] | Generic | n.s. | Model initiators, computing partners, validators | | ✓ | | | | | |
| [67] | Finance | CS | Follower candidates, leader nodes | ✓ | ✓ | ✓ | | ✓ | | |
| [68] | IoT | CD | Users, edge server, cloud server | | | | | | | |
| [69] | Generic | CD | Worker nodes | ✓ | ✓ | | | | | ✓ |
| [70] | Generic | n.s. | Task publishers, parties, miners, smart contract | ✓ | ✓ | | | | | ✓ |
| [71] | Generic | n.s. | Requesters, workers, crowdsourcing platform | | ✓ | | | ✓ | ✓ | ✓ |
| [72] | IoT | n.s. | IoT devices, edge servers, central cloud server | | ✓ | | | | | ✓ |
| [23], [24] | Generic | n.s. | Task requesters, workers | ✓ | ✓ | | | ✓ | | |
| [73] | Generic | n.s. | Requesters, workers | | ✓ | | | | ✓ | |
| [74] | IoV | CD | RSUs (aggregators), MECs, vehicles (workers) | | ✓ | | | ✓ | ✓ | |
| [75] | IoT | n.s. | Requesters, edge servers (workers), data collectors | | ✓ | | | ✓ | | ✓ |
| [76] | IoV | n.s. | UAVs (sensors, workers), requesters, MECs | | ✓ | ✓ | | | ✓ | ✓ |
| [77] | Generic | n.s | Requesters, data-arbitrators, workers | | ✓ | | | | | |
| [78] | Generic | n.s. | Requesters, miners | | ✓ | ✓ | | ✓ | ✓ | |
| [79] | Generic | n.s. | Miners, workers | | | ✓ | | | | |
| [80] | Generic | n.s. | Requesters, workers, aggregation servers | | ✓ | | | ✓ | ✓ | ✓ |
| [81] | Generic | n.s. | Administrator, requesters, workers, miners | | ✓ | | | ✓ | | |
| [82] | Generic | CS, CD | Workers, miners | | ✓ | | ✓ | | | |

models [59]. Kansra *et al.* integrate data augmentation, a technique to synthetically generate data such as images, into FL to increase model accuracy for ITS such as autonomous driving and road object identification [64]. Wang *et al.* propose a FLF dedicated to the crowdsensing of Unmanned Aerial Vehicles (UAVs) [76]. As UAVs are often equipped with multiple sensors and can be easily deployed to sensing areas, a FLF with UAVs has a huge benefit for ITS applications such as traffic monitoring and public surveillance.

Finance is the other domain that we found in the surveyed papers. He *et al.* proposed a FLF for commercial banks to better utilize customers' financial information [67]. Financial information such as credit level, risk appetite, solvency, movable and real estate owned are crucial sources to understanding the characteristics of customers of financial services. However, it is too sensitive to directly use them for data mining. Hence, a FLF is a viable framework for financial information management.

### 4.1.2 RQ 1-2: What problems were solved?

The problems solved by the papers can be categorized into (*i*) the SPoF issue in FL, (*ii*) BC-related issues, (*iii*) lack of clients' motivation, (*iv*) how to fairly evaluate clients' contribution, (*v*) security and privacy issues.

Most of the papers (29 out of 40) propose a system architecture of FLF to solve the problems of a SPoF in the current centralized server-clients architecture. More specifically, this issue is rooted in the structure of the original FL where an aggregation server is collecting local model updates from clients in a centralized manner. The idea to mitigate this issue is to decentralize the processes involved in FL using BC technologies. Each paper proposes operations, functions, and protocols processed in and outside the smart contract. Furthermore, some solve the issue of scalability in the FLF (e.g., [59], [63]) and BC-

related issues such as energy waste of consensus algorithms (e.g., [78], [79]). We will go into the proposed system architectures and BC-related issues in Section 4.2.

An incentive mechanism is integrated into FLFs to solve the problem of a lack of client motivation. The basic idea is to give monetary incentives to clients in return for their effort in training a local model. The incentive mechanism is also leveraged to solve the model poisoning attack which is an attack on a model update to deteriorate the quality of a global model by malicious clients' providing bogus local model updates. The idea for demotivating such attacks is to devise an incentive mechanism that penalizes malicious activities. Furthermore, a reputation score based on contribution is also useful to screen potentially malicious clients. Here, we need a contribution measurement metric to fairly evaluate the quality of clients' model updates and detect the attacks. Details about the proposed incentive mechanisms and contribution measurements will be covered in Section 4.3.

20 out of 40 papers propose approaches to solve issues related to security and privacy. With few exceptions (i.e., attacks on reputation [60], [65] and [79]), both security and privacy issues are rooted in local model updates. The security issue is related to the model poisoning which we mentioned above, while the privacy issue is related to sensitive information that might be leaked from the local updates. We will further summarize the works that solve the security and privacy issues in Section 4.4.

## 4.2  RQ 2: Blockchain

### 4.2.1   RQ 2-1: What is the underlying BC architecture?

Table 4 shows the overview of BC features. The BC system and its underlying consensus mechanism are an influential part of the FLFs infrastructure. FLFs are heterogeneous in terms of architecture, operation and storage requirements, contribution calculation, actors, and applied cryptography. Customized and tailored BC solutions may be required with respect to the underlying use case. Due to its restrictive scalability in terms of computation and storage, most of the analyzed FLFs apply BC as a complementary element in a more complex system, with a few exceptions [22], [23], [24], [55], [69]. BC systems themselves are complex distributed systems, heterogeneous across many dimensions, yet can roughly be categorized into public, private, and permissioned BCs.

1) **Public**: BCs are open access where participants can deploy contracts pseudo-anonymously
2) **Private**: BCs do not allow access for clients outside the private network and require an entity that controls who is permitted to participate
3) **Permissioned**: BCs are private BCs with a decentralized committee that controls the onboarding process

Note that the FLFs that utilize open-source public BCs such as Ethereum [50], [51], [53], [55], [56], [58], [65], [69], [71], [72], [77], Stellar [66] and EOS [73] were not deployed on the respective public BC in the experiments due to the enormous costs this would incur. Hyperledger Fabric [30] or Corda [84] are permissioned BCs running on

private networks, allowing for faster throughput through a limited amount of potential nodes. This makes these frameworks more suitable for applications where BC replaces computationally expensive operations such as aggregation or storage of neural network models.

The consensus protocol ensures the alignment and finality of a version across all distributed nodes without the need for a central coordination entity. While PoW is the most common mechanism applied in Bitcoin and Ethereum, it comes at the cost of wasting computational power on brute-forcing algorithmic hash calculations for the sole purpose of securing the network. Since many operations within the FLF frameworks are computationally expensive, these tasks can be integrated into the consensus mechanism which creates synergy and might be a better use of resources. Examples of consensus mechanisms can be found where the model accuracy is verified (Proof-of-Knowledge [59]), reputation scores are checked (Proof-of-Reputation [54]), the model parameters are securely verified (Proof-of-Federated-Learning [78]), the Shapley value is calculated for contribution measurement [25] or verification of capitalizing on efficient AI hardware (Proof-of-Model-Compression [79]).

### 4.2.2   RQ 2-2: How is BC applied within the FLF and what operations are performed?

BC technology is applied to mitigate the single point of failure and power imbalance of the server-worker topology of traditional FL through a transparent, immutable, and predictable distributed ledger. Embedded cryptocurrencies suit the useful property of real-time reward payments for predefined actions at the same time. In general, Turing-complete smart-contract-enabled BCs allow for a variety of possible complementary features for the FL training process, namely aggregation, payment, coordination, and storage:

1) **Aggregation**: The aggregation of model parameters, can be performed by a smart contract on top of BC [19], [21], [23], [24], [63], [67]. Since BC is assumed to be failure resistant, this strengthens the robustness against possible single-point of failure of an aggregation server. In addition, the deterministic and transparent rules of smart contracts ensure inherent trust with an equal power distribution among participants, while the transparency ensures auditability of contributions. Yet since every node in the BC has to compute and store all information, submitting a model to the smart contract for aggregation causes overhead in terms of both computation and storage on the BC. Assuming $n$ FL-workers and $m$ BC nodes over $t$ rounds, the BC scales with $\mathcal{O}(t * n * m)$ which questions the feasibility of on-chain aggregation.

   There are two papers that try to reduce data size for on-chain aggregation. Witt *et al.* [22] proposed a system where 1-bit compressed soft-logits are stored and aggregated on the BC saving communication, storage, and computation costs by orders of magnitude.

   Feng *et al.* [63] employ a framework based on two BC layers where the aggregation process is outsourced to a mobile edge server.

TABLE 4: Overview of BC features.

| Ref. | Operations | | | | BC | Consensus | On-chain | | Off-chain | Scalability |
|---|---|---|---|---|---|---|---|---|---|---|
| | Agg. | Cor. | Pay. | Str. | | | Items | Eval. | | |
| [22] | ✓ | ✓ | ✓ | | Agnostic | n.a. | 1-bit results of all participants | ✓ | | ✓ |
| [19] | ✓ | | ✓ | | Corda V3.0 | Algorand | Gradients | | | |
| [50] | | | ✓ | | Ethereum | PoW | n.s. | | ✓ | |
| [51] | | | ✓ | | Ethereum | PoW | Contribution, Merkle tree | | | |
| [52] | | ✓ | ✓ | | n.s. | n.s. | Models | | ✓ | |
| [53] | | | ✓ | | Ethereum | PoW | n.s. | ✓ | | |
| [54] | | | | ✓ | TrustRE | PoR | Reputation scores | | | |
| [20] | | | | ✓ | Custom | PoW | Model updates | | | |
| [25] | | | ✓ | | Custom | PoSap | SV values, tasks | | | |
| [55] | | ✓ | ✓ | ✓ | Ethereum | PoW | Tasks | | ✓ | |
| [21] | ✓ | | ✓ | ✓ | n.s. | PoW | Model updates, metadata | | | |
| [56] | | | ✓ | | Ethereum, HF | PoW | PoT records | ✓ | | |
| [57] | | ✓ | | ✓ | n.s. | n.s. | Reputation scores | | ✓ | |
| [58] | | ✓ | ✓ | | Ethereum | PoW | Users' addresses | | | |
| [59] | | ✓ | | ✓ | n.s. | PoK | Local models, loss, signatures | | | ✓ |
| [60] | | | ✓ | | Corda V4.0 | PBFT | Reputation scores | | ✓ | |
| [61] | | ✓ | ✓ | ✓ | n.s. | Custom | Client info, model parameters | | | ✓ |
| [62] | ✓ | | | | n.s. | n.s. | Masked gradients, global models | | | |
| [63] | ✓ | | ✓ | ✓ | HF | Raft | Local updates | | | ✓ |
| [64] | | ✓ | ✓ | | n.s. | n.s. | | | | |
| [65] | | ✓ | ✓ | | Ethereum | PoW | IPFS CIDs of models | | ✓ | |
| [66] | | ✓ | ✓ | | Stellar | n.s. | IPFS CIDs of models | | ✓ | |
| [67] | ✓ | ✓ | | ✓ | Custom | Raft | Local and global models, loss | | | |
| [68] | | | | ✓ | n.s. | Custom | Local and global models | | | |
| [69] | ✓ | ✓ | | | HF | Raft | Parameters | | | ✓ |
| [69] | | | ✓ | ✓ | Ethereum | PoW | Hashes of parameters | | | ✓ |
| [70] | | ✓ | ✓ | ✓ | Ethereum | PoW | Aggregated models | | | |
| [71] | | ✓ | ✓ | ✓ | Ethereum | PoW | Encrypted models | ✓ | | ✓ |
| [72] | | ✓ | | | Ethereum | PoW | Aggregated local updates | | | |
| [23], [24] | ✓ | ✓ | ✓ | ✓ | Agnostic | n.a. | Tasks, voting results, model updates | | ✓ | |
| [73] | | | ✓ | ✓ | EOS, HF | n.s. | Hashes of model updates, data size | | ✓ | |
| [74] | | ✓ | ✓ | ✓ | n.s. | n.s. | Models | | | |
| [75] | | | ✓ | | Agnostic | PBFT | Model updates | | | |
| [76] | | | ✓ | | n.s. | PoW | Tasks, model updates, aggregated models | | | |
| [77] | | | ✓ | | Ethereum | PoW | Reputation scores | | ✓ | |
| [78] | | ✓ | ✓ | | Custom | PoFL | Model updates | | | |
| [79] | | | ✓ | ✓ | Custom | PoMC | Model updates | | | |
| [80] | | | ✓ | | n.s. | n.s. | Signatures, reputation scores, contributions | | | |
| [81] | | ✓ | ✓ | ✓ | Agnostic | n.s. | Tasks, voting results, model updates | | | |
| [82] | ✓ | | ✓ | ✓ | Custom | PoW | Model updates, computation time | | | |

2) **Coordination**: Applying BC to coordinate and navigate the FL process allows for decentralization without the heavy on-chain overhead.

Instead of aggregating the model on-chain, letting the BC choose a leader randomly can ensure decentralization [52], [59], [69], [81]. Another way BC coordinates the FL process is by enabling the infrastructure for trustless voting atop the BC. Voting on the next leader (aggregator) [66], [67] or on each other's contributions [23], [24], [81] further democratizes the process. Beyond explicit coordination operations like voting or leader selection, the implicit function of storing crucial information and data for the FL process, [22], [55], [71], [74], verifying the correctness of updates [52], [69] or keeping the registry of active members [19], [22], [23], [24], [25], [55] is crucial for the FL workflow and implies coordination through BC as an always accessible, verifiable, transparent and immutable infrastructure.

3) **Payment**: Many general-purpose BC systems include cryptocurrency capabilities and therefore allow for the incorporation of instant, automatic and deterministic payment schemes defined by the smart contract's business logic. This advantage was capitalized on by 26 of the 40 FLF we analyzed.

Section 4.3 discusses the details of applied payment schemes in the context of reward mechanisms and game theory.

4) **Storage**: Decentralized and publicly verifiable storage on the BC facilitates auditability and trust among participants. Even though expansive, since all BC nodes store the same information in a redundant fashion, it might make sense to capitalize on the immutability and transparency feature of BC and store information where either shared access among participants is required or where verifiability of the history is required to hold agents accountable for posterior reward calculations [54]. In particular, machine learning models [19], [20], [21], [23], [24],

[59], [61], [62], [67], [68], [69], reputation scores [54], [57], User-information [19], [22], [23], [24], [25], [55] and Votes [22], [23], [24] are stored on-chain of the respective FLF.

### 4.2.3 RQ 2-3: Is the framework scalable?

Especially if the FLF is intended to be used with hundreds to millions of devices, the scalability of the framework is an important characteristic. In particular, (*i*) storing large amounts of data such as model parameters and (*ii*) running expansive computations on the BC e.g., aggregating millions of parameters, calculating expansive contribution measurements like Shapley Value or privacy-preserving methods hinder the framework to scaling beyond a small group of entrusted entities towards mass adoption. To overcome the scalability bottleneck of storage, some FLF applied an Interplanetary File System (IPFS) [85], where data is stored off-chain in a distributed file system, using the content address as a unique pointer to each file in a global namespace over computing devices [50], [52], [55], [57], [60], [65], [66], [73], [77]. Other FLFs are based on novel design choices to tackle the scalability issues: Witt *et al.* [22] applied compressed Federated Knowledge Distillation, storing only 1-bit compressed soft-logits on-chain. Chai *et al.* [59] design a hierarchical FLF with two BC layers to reduce the computational overhead by outsourcing computation and storage to an application-specific sub-chain. Similarly, Feng *et al.* [63] propose a two-layered design, where the transaction efficiency of the global chain is improved through sharding. Bao *et al.* [61] employ an adaption of Counting Bloom Filters (CBF) to speed up BC queries in the verification step of their FLF. Desai *et al.* [69] combine public and private BCs, with the former storing reputation scores for accountability and the latter used for heavy computation and storage. Furthermore, the authors apply parameter compression for further scalability improvements.

### 4.3 RQ 3: Incentive Mechanisms

#### 4.3.1 RQ 3-1: How are incentive mechanisms analyzed?

In general, the analysis comprises three steps: The first step is to determine what entities' behavior is examined. In FL, such entities could be workers or task requesters. The second step is to model the entities' utilities or profits. They can be obtained by taking the expectation of possible profits and costs into account. The last step is to analyze the defined utilities or profits. This could be done in a theoretical manner and/or via simulation. 30% (12 out of 40) of the surveyed papers analyzed the incentive mechanism theoretically, while 45% (18 out of 40 papers) measured workers' rewards via simulation. However, we only focus on the papers with theoretical analysis in this section as we do not see much technical depth or differences in the simulation-based analysis.

In the first step, 28 papers assume that only workers exist, while 12 FLFs have additional entities that pay rewards (e.g., task requesters) [23], [24], [54], [60], [61], [70], [73], [75], [77], [78], [80], [81]. In such a case, the utilities of both entities have to be analyzed such that task requesters can be profitable even if they pay rewards to workers. For instance, workers and task requesters are assumed in [23],

and it is vital to determine their behavioral assumptions (e.g., their goals, rationality, etc.).

The second step is to define the utilities or profits of entities. A utility is a one-dimensional measurable unit that quantifies an entity's value on an outcome and can have positive (e.g., rewards for workers and a value of AI models for task requesters) and negative values (e.g., a computation cost for workers and a total amount of payout for task requesters). Utilities and profits can be derived by subtracting costs from payouts (Eq. (2)). Although the elements of payouts $\Pi$ are mostly straightforward (e.g., rewards for work contribution), the cost elements $C$ are dependent on the assumed application scenarios. Typical costs in the surveyed papers are computation, electricity (e.g., [21], [23], [24], [59], [60], [76]), data acquisition (e.g., pictures and sensor readings [74], [75], [76]) and privacy leakage due to model updates (e.g., [75], [76], [78]). Even multiple cost factors can be considered (e.g., [74], [75], [76]).

The last step is to analyze the defined utilities and profits to ensure the robustness of the incentive mechanisms and derive the optimal reward allocation. A simple yet crucial analysis would be to prove that it is worthwhile for workers to join a FL task by showing that their profits are non-negative. For instance, Bao *et al.* modeled requesters' and/or workers' profits with given rewards and costs and proved that their profits are non-negative [61]. Utilities can be used to derive task requesters' and workers' optimal strategies by finding a point where utilities are maximized. By proving the existence of such a point, an equilibrium can be derived, which is a condition where entities (e.g., workers) cannot be better off deviating from their optimal strategies. An equilibrium state, if existent, is proof that a designed mechanism is stable. An equilibrium can be found by deriving the first- and second-order derivatives of the utility function with respect to the variable in question. For instance, Toyoda *et al.* optimizes the workers' data size [24].

Other works that determine optimal prices for tasks: Wang *et al.* propose a Q-learning-based approach to determine the optimal prices so that utilities are maximized via iterative learning processes [76]. Similarly, Zou *et al.* derive the optimal prices for workers with first- and second-order conditions when the value of data, transmission quality, and communication delay are the factors to determine their competitiveness and costs [74]. Hu *et al.* propose a two-stage optimization method to determine the optimal values of data and their prices in order by solving an Euler-Lagrange equation of their utilities. These kinds of two-stage optimization games are often formalized as a Stackelberg game [21], [59]. Jian and Wu propose a Stackelberg-game-based incentive mechanism for FL [21]. They analyzed the equilibria of two reward policies where a contribution is measured by data size or accuracy by modeling an aggregation server as leader and workers as followers. The uniqueness of their analysis is that they incorporate training and uploading time into the analysis in terms of a constraint as each round has a deadline in the FL.

Chai *et al.* propose a multi-leader multi-follower Stackelberg game to analyze their incentive mechanism in IoV [59]. Aggregation servers (RSUs) are leaders while workers (vehicles) are followers, and aggregation servers

TABLE 5: Overview of incentive mechanisms and contribution measurement.

| Ref. | Sim. | Theoretical analysis | Costs | Contribution Metrics | Abs. | Rel. | Rep. | Validator |
|------|------|----------------------|-------|---------|------|------|------|-----------|
| [22] | ✓ | ✓ | Generic | Correlation of predictions (Peer truth serum) | | ✓ | | Smart contract |
| [19] | | ✓ | n.s. | n.s. | | | | Miners |
| [50] | | | n.a. | Accuracy (validation scores) | ✓ | | | Miners |
| [51] | ✓ | | n.a. | Accuracy, data size | ✓ | | | Agg. server |
| [52] | ✓ | | n.a. | Euclidean distance of model updates | | ✓ | | Miners |
| [53] | ✓ | | n.a. | Data size | ✓ | | | Smart contract |
| [54] | ✓ | | n.a. | Accuracy, energy consumption, data size | | ✓ | ✓ | Task requesters |
| [20] | | | n.a. | Accuracy (loss) | ✓ | | | MEC servers |
| [25] | ✓ | | n.a. | Accuracy (loss) (Shapley values) | | ✓ | | Miners |
| [55] | | | n.a. | Accuracy (loss, marginal) (rank) | | ✓ | | Validators |
| [21] | ✓ | ✓(Stackelberg game) | Computation | Accuracy, data size | ✓ | | | Agg. servers |
| [56] | | | n.a. | n.s. | | | | n.s. |
| [57] | | | n.a. | n.s. | | | | n.s. |
| [58] | | | n.a. | Data size | ✓ | | | Task requesters |
| [59] | ✓ | ✓(Stackelberg game) | Computation | Accuracy (loss) | ✓ | | | Agg. servers |
| [60] | ✓ | ✓(Contract theory) | Energy | RONI [86], FoolsGold [87] | | | ✓ | Task requesters |
| [61] | | ✓ | Generic | n.s. | | | | Workers |
| [62] | ✓ | | n.a. | Generic (Shapley values) | | ✓ | | Smart contract |
| [63] | | | n.a. | n.s. | | | | n.s. |
| [64] | | | n.a. | Accuracy (generic, marginal) | ✓ | | | n.s. |
| [65] | ✓ | | n.a. | Accuracy (generic) | | ✓ | | Workers |
| [66] | | | n.a. | n.s. | | | | Validators |
| [67] | ✓ | | n.a. | Similarity of model updates (Shapley values) | | ✓ | | n.s. |
| [68] | ✓ | | n.a. | Communication delay, energy consumption | ✓ | | | Agg. servers |
| [69] | | | n.a. | Speed of model submission | ✓ | | | n.s. |
| [70] | | | n.a. | n.s. | | | | n.s. |
| [71] | | | n.a. | Accuracy, data size | ✓ | | ✓ | Task requesters |
| [72] | | | n.a. | Data size | ✓ | | | Agg. servers |
| [23], [24] | | ✓(Contest theory) | Computation | Accuracy (generic) (rank by voting) | | ✓ | | Workers |
| [73] | | | n.a. | Data size | ✓ | | | Workers |
| [74] | ✓ | ✓ | Data, communication | Accuracy (loss) | ✓ | | | Agg. servers |
| [75] | ✓ | ✓ | Sensing, privacy | n.s. | | | | n.s. |
| [76] | ✓ | ✓(Reinforcement learning) | Sensing, privacy, energy | Data size, sensing capacity | ✓ | | | Agg. servers |
| [77] | | | n.a. | n.s. | | | | n.s. |
| [78] | ✓ | ✓ | Privacy | Accuracy | ✓ | | | BC nodes |
| [79] | | | n.a. | n.s. | | | | n.s. |
| [80] | ✓ | ✓ | n.s. | Accuracy (loss) | ✓ | | ✓ | Workers |
| [81] | | | n.a. | Accuracy (generic) (rank by voting) | | ✓ | | Workers |
| [82] | | | n.a. | Computation time | ✓ | | | Miners |

first suggest prices and workers determine how much data they should collect and use for training so that both entities' utilities are maximized in order. Due to the high dimensionality of each worker's strategy, it is difficult to employ the traditional backward induction method to derive an equilibrium. Hence, they leverage the Alternating Direction Method of Multipliers (ADMM) algorithm [88] to iteratively reach the social optimum point.

Three papers model the incentive mechanism in FL as a contract or contest: a task requester proposes a contract with a task description and its reward and workers can determine whether or not to sign such a contract and how many resources they will provide [60]. A FL process can be also seen as a contest as workers need to work first, which incurs irreversible costs due to computation, whereas their rewards are not guaranteed at the time of model update submission. Toyoda *et al.* give an incentive analysis based on the contest theory [23], [24]. Workers' utilities are used to derive how much effort workers' should exert on a task under the risk of not gaining prizes, while requesters' utility is used to determine how a prize should be split among workers.

### 4.3.2 RQ 3-2: How are the contributions of clients measured?

Incentive mechanisms require (*i*) the measurement of contribution by each client to (*ii*) fairly distribute rewards. However, the clients' contributions in form of model updates or gradients do not imply direct information on the overall performance metric like the accuracy of the global model.

The metrics used in the literature can be categorized into absolute and relative ones. The absolute metrics are metrics that can be measured without others' local model updates. For instance, a loss function can be measured from a local model and a global model, and the difference between them can be used as a metric for contribution measurement. Although the majority of absolute metrics are based on the accuracy (e.g., [50], [51]) and data size (e.g., [21], [53]), other factors are also proposed such as energy consumption (e.g., [54], [68]) and computation time [82]. Some combine multiple metrics (e.g., [54], [71]). Absolute metrics are generally straightforward but hard to validate, e.g., metrics that are based on the data size used depending on the client's honesty. In contrast, relative metrics can be measured by comparing submitted results

(e.g., gradients, model updates) in terms of correlation or ranking. For instance, Zhao *et al.* propose a metric based on the Euclidean distance of workers' model updates [52]. Likewise, Witt *et al.* utilizes peer-truth serum [89] in the FL context, where contributions are measured based on the correlation of prediction on the labels of a public dataset [22]. Voting is another approach to determine clients' contribution relatively. For instance, clients choose the best model updates from the previous `FedAvg` round by ranking the respective updates based on the accuracy using their local datasets [23], [24], [81].

A similar metric is clients' reputations. If the same clients are assumed to join different FL tasks, reputation scores calculated based on clients' past contributions can be used to determine the reward distribution (e.g., [52], [54], [60], [71]). For instance, Kang *et al.* propose to calculate workers' reputation based on a direct opinion by a task requester and indirect opinions by other task requesters [60].

However, even if the clients' individual contribution is measured, the question regarding fair distribution remains an unsolved issue. The Shapley value is an approach to determine payouts to workers based on the marginal utility added, taking all possible combinations of contributors into consideration [90]. Three papers propose to use the Shapley value for fair reward distribution in the FL context [25], [62], [67]. Liu *et al.* applies the Shapley value based on the accuracy of a test dataset [25]. He *et al.* compared their Shapley-value-based method with three approaches, namely (*i*) equal distribution, (*ii*) a method based on individual contribution, and (*iii*) a method called the labor union game where only the order of submission is taken into account to contribution measurement, and found that the Shapley-value-based method outperforms the others in terms of workers' motivation and fairness [67]. Ma *et al.* propose a method to calculate Shapley values even if model updates are masked to preserve workers' privacy [62].

Which entities validate the contribution is an open issue, complementary to the contribution measurement. The validators in FLF can be classified into (*i*) aggregation servers FL (e.g. [20], [51], [59]), (*ii*) task requesters (e.g. [54], [60], [71]), (*iii*) validators whose task is only to measure contribution ( [55], [66]), (*iv*) BC nodes (e.g., [50], [78], [82]), (*v*) workers (e.g., [65], [80], [81]) and (*vi*) smart contracts (e.g., [22], [53], [62]). Some of the works assume that aggregation servers, task requesters, or validators are expected to possess datasets to calculate the metrics discussed above. As reviewed in Section 4.2, others propose custom BC architectures for FL where the validation process is integrated into the consensus mechanism, making BC nodes validators. In some scenarios, aggregation servers, task requesters, and BC nodes take up the role of validators since they aggregate the model updates. However, datasets for validation may not be always available. Furthermore, metrics based on the correlation of predicted labels do not require any validation dataset and can even be measured in a smart contract [22].

## 4.4 RQ 4: Experiments

Conducting experiments is a key element of FLF development for two reasons. Firstly, the implementation of an example testifies to the feasibility of the approach and gives the authors the chance to identify weaknesses of their frameworks, e.g., poor scalability. Secondly, conducting experiments allows the comparison of the proposed approaches with each other, e.g., based on the accuracy of the models on standardized test sets. We screened the papers for experiments, and when present, examined them according to nine criteria (TABLE 6).

### 4.4.1 RQ 4-1: Is the performance of the framework reported?

The large majority of papers report results of their experiments expressed in either loss, accuracy, or F1 score (84.8%, 28 out of 33 papers with experiments). The remaining instead focus on the performance of their novel group-based Shapley value calculation for contribution measurement [62], the user interface [25], the computational effort and adversarial influence [69], or game-theoretic quantities such as utility values and rewards [76]. We note that comparability of the approaches is not given through the conducted experiments, since even when using the same data sets and the same evaluation metrics, different experimental scenarios are investigated. In conclusion, to obtain insightful results, experiments should compare the performance (accuracy and computational effort) of an incentivized, decentralized FL system in a standardized challenging environment (non-IID, adversaries) with either the performance of a traditional centralized FL system or the performance of a locally trained model without FL. Ideally, the effectiveness of IM and decentralization efforts are reflected in the FLF performance through a holistic experimental design.

### 4.4.2 RQ 4-2: How comprehensive are the experiments?

First, it was found that the majority of publications do include experiments. Only seven of 40 papers did not conduct experiments [23], [24], [57], [64], [66], [67], [77]. However, the analysis also shows that only 45% of the experiments implement the actual BC processes (15 out of 33 papers with experiments). Instead, the distributed functionality was simulated or its impact estimated. For instance, Mugunthan *et al.* [65] focus on the evaluation of the frameworks' contribution scoring procedure by simulating collusion attacks on the FL procedure. The effect of introducing BC to the FL framework was accounted for by estimating the per-agent gas consumption. Similarly, Chai *et al.* [59] conducted experiments specifically designed to investigate the Stackelberg game-based incentive mechanism. The authors accomplish this without implementing the BC processes.

To test the FL functionality of the framework, an ML problem and a dataset must be selected. For the ML application and the dataset used, we observe a high homogeneity. Almost all experiments realize classification problems and use publicly available benchmark datasets. The most common are MNIST (handwritten digits) [92] and its variations, as well as CIFAR-10 (objects and animals) [93]. Only Rathore *et al.* [72] and Li *et al.* [20] did not perform classification tasks. Rathore *et al.* [72] performed object detection on the PASCAL VOC 12 dataset [94]. Object detection typically combines regression and classification by

TABLE 6: Overview of experiments. *(BCWD = Breast Cancer Wisconsin Data Set, BT = Blockchain Tampering, DGHV = Dijk-Gentry-Halevi-Vaikutanathan Algorithm, DP = Differential Privacy, ECC = Elliptic Curve Cryptography, HE = Homomorphic Encryption, HDD = Heart Disease Data Set, KDD = Knowledge Discovery and Data Mining Tools Competition, RP = Random Model Poisoning, RSA = Rivest-Shamir-Adleman Cryptosystem, RT = Reputation Tampering, SA = Secure Aggregation [91], SP = Systematic Model Poisoning, ZKP = Zero-Knowledge Proof, 2PC = 2-Party Computation)*

| Ref. | Tasks | | Datasets | #Clients | Algorithms | Privacy | Non-iid | Adversaries | | | | Imp. | Per. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Clf. | Rgr. | | | | | | BT | RP | RT | SP | | |
| [22] | ✓ | | EMNIST | 10 | FD | | ✓ | | ✓ | | ✓ | | ✓ |
| [19] | ✓ | | MNIST | 4-10 | FedAvg | HE (Paillier) | | | | | | ✓ | ✓ |
| [50] | ✓ | | MNIST | 5 | FedAvg, EWC | DP, HE | ✓ | | | | | | ✓ |
| [51] | ✓ | | Original | 4 | FedAvg, CDW | n.s. | ✓ | | | | | ✓ | ✓ |
| [52] | ✓ | | MNIST | 10 | FedAvg | DP | | | ✓ | | | | ✓ |
| [53] | ✓ | | MNIST | 25 | FedAvg | n.s. | | | | | | ✓ | ✓ |
| [54] | ✓ | | MNIST, CIFAR10 | n.s. | n.s. | | | | | | | | ✓ |
| [20] | | ✓ | Real-time AD video | 1 | n.s. | HE (DGHV), ZKP | | | | | | | ✓ |
| [25] | ✓ | | MNIST | n.s. | FedAvg | SA | | | | | | ✓ | |
| [55] | ✓ | | MNIST | 5 | FedAvg, FedProx | Sym. cryptography | | | ✓ | | | ✓ | |
| [21] | ✓ | | Reddit, Celeba | 5-75 | n.s. | | | | | | | | ✓ |
| [56] | ✓ | | MNIST | 100 | Original | Ring sig., HE (RSA), Rabin, ECC | | | | | | ✓ | ✓ |
| [57] | | | n.a. | n.a. | n.a. | | | | | | | | |
| [58] | ✓ | | Mathworks handwritten | 10 | FedAvg | Pairing-based cryptography, ECC | | | | | | ✓ | ✓ |
| [59] | ✓ | | MNIST, CIFAR10 | 6 | Original | Asym. cryptography, signatures | | | ✓ | | | | ✓ |
| [60] | ✓ | | MNIST | 100 | FedAvg | | ✓ | | ✓ | ✓ | | ✓ | ✓ |
| [61] | ✓ | | n.s. | 10 | n.s. | n.s. | | | | | | ✓ | ✓ |
| [62] | ✓ | | ORHD | 9 | FedAvg | SA | ✓ | | | | | | |
| [63] | ✓ | | MNIST | 30 | n.s. | | | | | | | | ✓ |
| [64] | | | n.a. | n.a. | n.a. | | | | | | | | |
| [65] | ✓ | | Adult census, KDD | 50 | Custom FedAvg | DP | ✓ | | ✓ | ✓ | ✓ | | ✓ |
| [66] | | | n.a. | n.a. | n.a. | | | | | | | | |
| [67] | | | n.a. | n.a. | n.a. | | | | | | | | |
| [68] | ✓ | | Original | 3, 4 | FedAvg | n.s. | | | ✓ | | | | ✓ |
| [69] | ✓ | | CIFAR10 | 100 | FedAvg, signSGD | n.s. | ✓ | | | ✓ | ✓ | | |
| [70] | ✓ | | FMNIST | 900 | FedAvg | HE (Paillier) | ✓ | | ✓ | | | ✓ | ✓ |
| [71] | ✓ | | BCWD, HDD | 3 | FedAvg | DP | | | | | | ✓ | ✓ |
| [72] | ✓ | ✓ | PASCAL VOC 2012 | 5-10 | FedAvg | HE | | | | | | ✓ | ✓ |
| [23], [24] | | | n.a. | n.a. | n.a. | | | | | | | | |
| [73] | ✓ | | MNIST | 10 | FedAvg | | ✓ | | | | | ✓ | ✓ |
| [74] | ✓ | | MNIST | 50 | n.s. | | | | | | | | ✓ |
| [75] | ✓ | | FEMNIST | 35, 105, 175 | n.s. | | ✓ | | | | | | ✓ |
| [76] | ✓ | | MNIST | n.s.[2] | n.s. | | | | | | | | |
| [77] | | | n.a. | n.a. | n.a. | | | | | | | | |
| [78] | ✓ | | CIFAR-10 | 20, 50, 100 | n.s. | HE, 2PC | | | | | | | ✓ |
| [79] | ✓ | | ImageNet | 20 | n.s. | | | ✓ | ✓ | | | ✓ | ✓ |
| [80] | ✓ | | MNIST, CIFAR-10 | 10 | FedAvg | | | | ✓ | | ✓ | | ✓ |
| [81] | ✓ | | MNIST | 50 | FedAvg | | | | | | | | ✓ |
| [82] | ✓ | | MNIST, CIFAR-10 | 2-6 | FedAvg | | ✓ | | | | | | ✓ |

predicting bounding boxes and labeling them. Li *et al.* [20] applied their FLF to autonomous driving and minimized the deviations in steering-wheel rotation between a human-driven and simulation-driven vehicle. This corresponds to a regression task.

As to the number of training data holders, the experiments considered between one [20] and 900 [70] clients. In general, one would expect papers specifying cross-silo settings to test with fewer (<100 [8]) and papers specifying cross-device settings to test with more (>100) clients. Of the frameworks clearly designed for a cross-device application, it is noticeable that only Kang *et al.* [60] and Desai *et al.* [69] conduct experiments with 100 participants or more. On the contrary, Rahmadika *et al.* [56] test with as many as 100 participants, although only designing a cross-silo framework.

Regarding the FL algorithm, the classic FedAvg [6] is mainly used. Furthermore, in some experiments algorithms are used that mitigate the problem of catastrophic forgetting (Elastic Weight Consolidation (EWC) [50]), reduce the communication overhead (Federated Knowledge Distillation (FD) [22], signSGD [69]), or show more robust convergence for non-IID and other heterogeneous scenarios

(FedProx [55], Centroid Distance Weighted Federated Averaging (CDW_FedAvg) [51]). Chai *et al.* [59] and Mugunthan *et al.* [65] design custom FL algorithms. Specifically, Chai *et al.* [59] propose a FLF with two aggregation layers in order to promote scalability. In their FL algorithm, nodes in the middle layer aggregate the local model updates of associated nodes in the lowest layer. This semi-global model is then fine-tuned by the middle layer nodes based on data collected by the middle layer nodes themselves. Finally, nodes in the top layer aggregate the fine-tuned models from the middle layer nodes into a global model which is eventually passed back to the lowest layer nodes. All aggregations are weighted by the training dataset size. Mugunthan *et al.* [65] propose a FLF where all clients evaluate and score the differentially encrypted locally trained models of all other clients. These scores are reported to a smart contract which computes an overall score for each local model. Eventually, each client aggregates the global model from all local models, weighted by the overall score.

### 4.4.3 RQ 4-3: Are non-IID scenarios simulated?

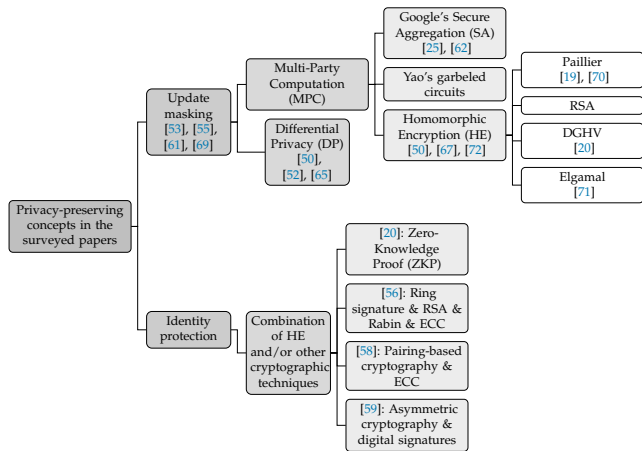In real-world applications of FL, the training data is often not Independent and Identically Distributed (non-IID)

Fig. 3: Privacy-preserving concepts employed in survey papers. *(DGHV = Dijk-Gentry-Halevi-Vaikutanathan Algorithm, ECC = Elliptic Curve Cryptography, RSA = Rivest-Shamir-Adleman Cryptosystem. Papers with unspecified methods are added to next highest node.)*

between the clients. This affects the performance of the global model and adds an additional layer of complexity with respect to contribution measurement. Hence, how we simulate non-IID scenarios with open datasets is crucial. In 11 publications and for various benchmark datasets, non-IID scenarios were considered. For the example of the MNIST dataset, Witt *et al.* [22] simulate different levels of non-IID scenarios following the Dirichlet distribution as it can easily model the skewness of data distribution by varying a single parameter. Martinez *et al.* [73] split the dataset in overlapping fractions of various sizes, whereas Kumar *et al.* [50] divide the dataset so that each trainer only possesses data from two of the ten classes. In a less skewed setting, Kumar *et al.* allocate data from at most four classes to each trainer, with each class being possessed by two devices.

### 4.4.4 RQ 4-4: Are additional privacy methods applied?
Even though FL's core objective is to maintain confidentiality through a privacy-by-design approach where model parameters are aggregated instead of training data, there remain innumerable attack surfaces [95]. Therefore, the presented frameworks employ additional privacy-preserving mechanisms which can be divided into two groups. (*i*) Mechanisms that encrypt or obfuscate gradients and prevent malicious parties to draw conclusions about the data set. (*ii*) Mechanisms that hide the identity of participating parties. A classification of the employed privacy-preserving methods can be seen in Fig. 3.

The methods of the first group can be further divided into (*i*) approaches that are based on cryptographic secure MPC, and (*ii*) approaches that are based on Differential Privacy (DP).

MPC refers to cryptographic methods by which multiple participants can jointly compute a function without having to reveal their respective input values to the other participants. MPC approaches include three groups of methods [91]. (*i*) Google's Secure Aggregation (SA) [91] is specifically designed to achieve low communication and computation overhead and to be robust towards device

dropout. It has been employed by Liu *et al.* [25] and Ma *et al.* [62]. Beyond implementing SA, the latter develops a group-based Shapley-value method for contribution measurement, since the native Shapley-value method cannot be applied to masked gradients. (*ii*) Yao's garbled circuits have not been applied to any of the analyzed frameworks, but are mentioned here for completeness. (*iii*) While Yao's garbled circuits were developed for 2-party secure computing, Homomorphic Encryption (HE) allows for higher numbers of participants [91]. As shown in Fig. 3, HE has been employed in several works [19], [20], [67], [71], [70], [72]. For that, different implementations of the homomorphic idea have been chosen, such as the Pallier cryptosystem, the Elgamal cryptosystem, or the Dijk-Gentry-Halevu-Vaikutanathan Algorithm (DGHV). Li *et al.* [71] choose the Elgamal cryptosystem that is less computationally expensive than other HE approaches.

DP refers to a method where noise is drawn from a probability density function $p_{noise}(x)$ with expected value $\mathbb{E}(p_{noise}(x)) = 0$ obfuscates the individual contribution with minimal distortion of the aggregation. DP has been employed by Mugunthan *et al.* [65], Zhao *et al.* [52], and Kumar *et al.* [50], with the latter combining the use of HE and DP. However, the fewer clients participate in the DP process, the heavier the distortion of the aggregated model, introducing a trade-off between privacy and model accuracy. Zhao *et al.* [52] mitigate the loss in accuracy by incorporating a novel normalization technique into their neural networks instead of using traditional batch normalization (e.g., [50], [65]). Besides MPC and DP, another technique for data set protection is chosen by Qu *et al.* [78]. Instead of the clients sharing masked gradients, the FLF relies on requesters sharing masked datasets in the model verification step. This prevents other workers from copying the models while testing and evaluating them. HE and 2-Party Computation (2PC) are used. Zhang *et al.* [55], Desai *et al.* [69], Bao *et al.* [61], and Rahmadika *et al.* [53] also rely on the masking of gradients but do not specify the privacy-preserving mechanisms.

The second group of frameworks targets the protection of participants' identities through cryptographic mechanisms. For that, Rahmadika *et al.* [56] combine ring signatures, HE (RSA), Rabin algorithm, and Elliptic Curve Cryptography (ECC), while Chai *et al.* [59] incorporate digital signatures and asymmetric cryptography approaches, and Rahmadika *et al.* [58] perform authentication tasks through pairing-based cryptography and ECC. Only one framework implements measures for both masking gradients as well as hiding identities. Li *et al.* [20] use DGHV for masking gradients and Zero-Knowledge Proof (ZKP) for identity protection.

Finally, He *et al.* [67] specifically address the problem of aligning entities. This problem occurs in vertical federated learning where different parties hold complementary information about the same user. The parties have to find a way of matching this information without disclosing the identity of their users. To solve this problem, He *et al.* employ Encrypted Entity Alignment which is a protocol for privacy-preserving inter-database operations [67].

### 4.4.5 RQ 4-5: Is the framework robust against malicious participants?

The experiments consider and simulate different types of adversaries, whereas some publications consider multiple types of attacks. Four groups of attack patterns were identified in the publications: random model poisoning, systematic model poisoning, reputation tampering (RT), and BC tampering. The most common attack considered in the experiments is random model poisoning. This includes attacks, where local models are trained on a randomly manipulated data set ( [55], [60], [65], [80]) or where random parameter updates are reported ( [68], [22], [60], [52], [59], [70], [79]). For instance, malicious agents in [55] use a training dataset with intentionally shuffled labels, whereas in [70] the parameter updates are randomly perturbed with Gaussian noise. Kang *et al.* [60] analyze the effects of a bad or manipulated data set by providing 8% of the workers with training data where only a few classes are present, and another 2% of the workers with mislabeled data. Kang *et al.* quantify the insufficiency of the dataset using the earth mover's distance.

The second most commonly simulated type of attack is systematic model poisoning where the attackers manipulate the model through well-planned misbehavior. In [69], a fraction of workers collude and manipulates their image classification data sets by introducing a so-called trojan pattern: the malicious agent introduces a white cross to a certain fraction of a class, e.g., to 50% of all dog pictures in an animal classification task and re-labels these data points as horse pictures. This creates a backdoor in the model that cannot be detected by subjecting the model to dog or horse pictures which will be correctly classified. However, pictures with the trojan pattern will be misclassified. Other forms of systematic model poisoning can be found with Witt *et al.* [22], Mugunthan *et al.* [65], Gao *et al.* [80].

The third type of attack that was simulated is Reputation Tampering (RT). Here, malicious agents intentionally provide colluding agents with perfect reputation or voting scores [60], [65]. The fourth type of attack is BC tampering [79]. Here, malicious miners intentionally fork the BC and prevail by building a longer branch faster than the honest miners.

### 4.5 RQ 5: Summary: What are Lessons Learned?

The inherent complexity of FLF leads to heterogeneity of the scientific research across the dimensions (*i*) application, (*ii*) overall design, (*iii*) special focus on open issues, and (*iv*) details and thoroughness.

**Application**: Although the majority of analyzed works offer application independent frameworks (classified as "generic" in Table 3) other FLFs are applied across IoT, Industrial-IoT (IIoT), IoV, and Finance. The heterogeneity of the required properties across those domains causes differences in the design choices of function, operations, storage of BC, contribution measurement, and privacy requirements.

**Variety of possible design choices**: In addition to the domain-specific influence on the system architecture, design choices about the FL algorithms, communication protocol, applications of BC within the ecosystem, BC technology (existing or novel), storage and operation on BC, security trade-offs, mechanism design, contribution measurement, etc. add to the complexity and overall variety of such systems. For example, some works apply BC as the outer complementary layer [54] while BC is the core infrastructure for coordination, storage, aggregation, and payment in other FLFs [22], [55]. Furthermore, some works developed application-specific BC systems, while others tried to embed a FLF on top of existing BC frameworks such as Ethereum for cheaper and pragmatic deployment. Our survey exposes a similar variety in the choice of the contribution measurement: The spectrum reaches from the computationally lightweight correlation of answers on a public dataset [22] as a proxy for contribution as opposed to the Shapley value, a measurement with strong theoretical properties but massive computational overhead [25], [67].

**Special Focus**: The aforementioned complexity as well as its novelty results in many open issues across a broad spectrum. Many works, therefore, focus on solving specific issues such as enhanced privacy [19], [20], [56], [66], novel BC systems [25], [51], bandwidth reduction [68], novel contribution measurements [22], [25], [68] or game theory (e.g., [21], [60]), as the major contribution which further complicates a holistic comparison of FLF.

**Thoroughness**: The analyzed papers also vary heavily in provided detail and thoroughness, ranging from first concepts, lacking details in terms of important specifications such as performance, specific function, operation and storage on BC, contribution measurement, robustness, experiments and privacy to theoretically detailed and experimentally tested solutions. None of the analyzed papers are production-ready.

### 4.5.1 Standards for better comparability

For better reproducibility, implementability, and comparability we suggest considering and defining the following elements when designing a FLF.

*System model and architecture:*

- Assumed application
- Type of FL (i.e., CD vs CS, horizontal vs vertical)
- Entities (including attackers)
- Setup (e.g., who manages a system, who deploys it)
- Role of BC within the FLF (e.g., what part does BC replace, what functions/operations)
- BC design (e.g., consensus algorithms, BCs, smart contracts)
- Non-BC design (e.g., off-chain storage, privacy protection, authentication)
- Procedures (e.g., flowcharts and diagrams)
- Theoretical analysis of incentive mechanisms
- Specification of clients' contribution measurement
- Possible attacks (e.g., system security, data privacy)

*Performance analysis:*

- Quantitative performance analysis
- Scalability analysis with respect to blockchain and contribution measurement

*Cost analysis:*

- Overhead and cost analysis of BC infrastructure
- Overhead and cost analysis of the contribution measurement
- Performance-cost trade-off discussion

## 5 FUTURE RESEARCH DIRECTIONS

The multitude of possible applications of FLF come with different requirements in terms of accuracy, latency, cost, and privacy. To account for this, we classify future research into two main directions, namely (i) increase in framework performance and (ii) expansion of framework functionalities.

### 5.1 Performance

Most state-of-the-art publications only consider BC and incentive mechanisms on a conceptual or theoretical level, however, they lack a performance analysis. Yet, low operational costs and latency, as required by real-time applications, such as autonomous driving, demand high-performance systems. To develop such frameworks, we have identified four performance bottlenecks as future research directions: (i) framework scalability, (ii) communications and network, (iii) framework implementation, and (iv) framework evaluation and comparison.

#### 5.1.1 Framework scalability

One of the major factors for the applicability of a FLF is its ability to scale beyond small groups toward mass adoption. Out of the 40 papers, only six mentioned and considered scalability within the design of their respective FLF. In particular, our reviews show that the integration of distributed ledger technology frequently leads to scalability problems. In FLFs, BC technology becomes a scalability bottleneck if

1) it is part of the operating core infrastructure of the FLF (e.g., [22]) and not only a complementary outer layer technology (e.g., [54])
2) heavy operations such as aggregation or reward calculation are performed on-chain [25]
3) a large amount of information is stored on the BC such as model updates
4) the BC framework is public and used outside the realm of the FLF
5) the consensus mechanism is resource-intense (e.g., PoW).

There are multiple promising future strategies to improve the scalability of the framework. First, FLF-specific BC systems have been proposed that replace the computational overhead of the PoW-based systems with computational heavy tasks in FLF such as model parameter verification [78], reputation verification [54], or contribution measurement calculations [25]. Secondly, Zhang et al. [79] have investigated the use of efficient AI hardware to increase BC scalability. Wang et al. explored the domain of resource optimization in BC-based FLFs to further improve the scalability [96]. Moreover, Weng et al. [19] aim to improve scalability by enhancing the privacy procedures

for the FLF processes. Another promising research direction is the application of Zero-Knowledge Succinct Non-Interactive Argument of Knowledge (ZK-SNARKs) [97] in the FLF context. ZK-SNARKs is a promising cryptographic technology that allows a *prover* to prove to a *verifier* that computation has been executed without revealing the program itself. This verification is faster than actually computing the original code and can be implemented easily on the smart contract. Hence, this will improve the scalability of BC-enabled FLF dramatically. However, due to its generality, which processes leverage ZK-SNRAKs is an open question. Finally, the performance of the BC itself can be improved, e.g., by increasing the number of transactions per second.

#### 5.1.2 Communication and network

Another major remaining challenge is the communication bottleneck. Decentralized wireless gadgets, as employed in decentralized FL, operate on lower communication rates than traditional intra- or inter-datacenter links. This leads to a trade-off between accuracy and communication cost. Although there exists first theoretical research on the nature of this trade-off, its findings have not yet been incorporated in the proposed FL frameworks [7]. In terms of communication rates, new developments are also expected once 6G technology is introduced, which is predicted to be mutually empowering with FL [83].

Furthermore, researchers face the communication-related problem of scheduling and resource allocation under dynamic channel condition and heterogeneous computing capacity of devices in IoT [98]. For instance, Yang et al. propose a device selection strategy in UAV to keep the low-quality devices from affecting the learning efficiency and accuracy [98].

Another challenge is related to key collisions during update communication: To avoid throughput issues, data is typically uploaded iteratively in multiple smaller batches, causing latency and collision effects to become more dominant. For instance, Desai et al. [69] point out that Hyperledger cannot deal with the Multiversion Concurrency Control (MVCC) of its underlying database so many transactions fail and need to be repeated. Accordingly, future research should be directed to the three compression objectives identified by Kairouz et al. [7]: gradient compression (client-to-aggregator communication), model compression (aggregator-to-client communication), and local computation reduction. In consequence, security and privacy mechanisms need to be adapted to operate on the compressed data (Section 4.4.4). Starting points for this upcoming research include sparsification and quantization approaches [99], 1-bit compression [22], or the parallelization on multiple contracts [59], [63], [69]. For the latter, Desai et al. [69] analyze the trade-off between communication speed and the number of employed parallelized contracts.

In general, future framework proposals should consider communication cost and time in their simulative methods, e.g., building up on Kang et al. [60].

### 5.1.3 Framework implementation

Most of the papers we reviewed are focused on the algorithm side. However, in order to go beyond theory towards real production-ready deployments, implementation details have to be taken into consideration. For instance, incurred deployment and maintenance costs of unproven novel BC systems are often ignored. Introducing a new, custom-made, and highly complex infrastructure introduces security risks and it requires a large team of experts to run and maintain such a system in practice. Therefore, software/hardware co-design is another vital topic in FL (e.g., [79], [100], [101], [102]). For instance, Wang *et al.* point out that cipher-text operation and encryption parts are major bottlenecks on the FL and proposed a novel Field Programmable Gate Array (FPGA) design for it [103]. We believe that there are potential research topics in the software/hardware co-design for FL. Interested readers may refer to the survey papers of Khan *et al.* [100], [101].

### 5.1.4 Framework evaluation and comparison

While many papers have conducted performance evaluation, few showed a comparison with other FLFs. This hinders the scientific advancement towards high-performing frameworks as the different design choices of the papers remain uncompared. Furthermore, the frameworks have often not been evaluated in realistic scenarios: (*i*) relatively well-known benchmark datasets such as MNIST and CIFAR-10 are chosen (29 out of 34 papers that conducted experiments on classification) and (*ii*) the non-IID setting is only applied in 11 out of 34 papers. Furthermore, inconsistencies between the targeted FL setting (i.e., CS, CD) and the number of clients in the experiments are observed. In particular, FLFs that assume CD should simulate a large number of clients, however, only Kang *et al.* [60] and Desai *et al.* [69] conducted experiments with 100 participants or more (TABLE 6).

To better evaluate FLFs, we suggest using common datasets dedicated to FL (e.g., LEAF [104]) as well as simulating different levels of non-IID data among clients (e.g., Dirichlet distribution [22]). We also suggest deploying a FLF on the clusters of inexpensive computers such as Raspberry Pi [105] to realistically simulate large-scale FL scenarios under the CD assumption.

Moreover, it is difficult to simulate the effect of decentralization and incentivization (e.g., Shapley value and game-theoretic mechanisms) in a comparable way since each paper uses different assumptions. Therefore, to fairly compare FLFs, holistic experiments should be designed, where the effects of decentralization and incentivization are captured by metrics such as overall accuracy, cost, or latency.

## 5.2 Functionalities

Future research should also focus on integrating further functionalities into the FLFs. Firstly, most of the proposed FL systems are limited to supervised classification, however, other types of ML problems should be considered as well. Secondly, lightweight privacy-preserving techniques are necessary for some applications that use sensitive information (e.g., medical logs and personal financial information). Thirdly, a fair, non-manipulable, and lightweight mechanism for contribution measurement has yet to be developed.

### 5.2.1 Beyond supervised FL and federated averaging

To expand the applicability of FLFs, machine learning tasks beyond supervised learning should be enabled, such as anomaly detection, reinforcement learning, natural language processing, user behavior analysis, and unsupervised learning tasks (e.g., [106], [107]). This will require new or adapted model aggregation algorithms and a new contribution measurement to integrate such tasks with IM and BC.

So far, `FedAvg` requires the same neural network architecture on all devices to participate. This may lead to issues in real-world environments where clients might have different hardware and bandwidth capabilities. Federated Knowledge Distillation [22] is an interesting novel FL approach in this context, allowing for a flexible neural network architecture and a dramatic reduction in bandwidth [108]. However, Federated Knowledge Distillation requires a public dataset to distill the knowledge.

Traditional deep learning algorithms such as DNNs and Convolutional Neural Networks (CNNs), are generally power-hungry, which is problematic in IoT environment. To address this challenge, biological neurons-inspired DNNs called Spiking Neural Networks (SNNs) have been actively studied for edge AI (e.g., [109], [110]). SNNs will enable edge devices to exploit brain-like biophysiological structure to collaboratively train a global model while helping preserve privacy. For instance, Lead Federated Neuromorphic Learning (LFNL) is a method to enable SNNs in a federated manner [110]. Furthermore, a leader election scheme is proposed to elect one device with high capability (e.g., computation and communication capabilities) as a leader to manage model aggregation, eliminating a fixed central coordinator and avoiding model poisoning attacks.

### 5.2.2 Towards lightweight privacy-preserving FL

Despite FL being a data privacy-preserving technology by design, research has shown that certain characteristics of the underlying training data sets can be inferred from the global model and that additional privacy-preserving measures are recommended. Our review shows that two classes of security concerns are targeted by the publications, namely (*i*) leakage of data set characteristics and (*ii*) disclosure of participant identities. Although a substantial number of papers (20 out of 40 publications) address one of these concerns, only a single paper addresses both [20]. Moreover, preventing data set leakage through DP or MPC inflicts trade-offs. Specifically, DP comes with a trade-off between data security and model accuracy, while MPC comes with a trade-off between data security and computation complexity, and it might thus not be applicable with a large number of participants [62]. It is worth noting that the model accuracy of DP cannot be inherently improved due to intentionally added noise. Hence, it would be important to explore lightweight MPC algorithms [91] to accommodate a large number of clients for privacy-preserving FL.

### 5.2.3 Towards fair, non-manipulable, and lightweight contribution measurement

Although multiple approaches for contribution measurements have been explored in the literature (Section 4.3), a fair, non-manipulable, and lightweight mechanism has yet to be developed as the following overview shows.

Firstly, a contribution can be measured based on the clients' honest reports of the amount of data, local accuracy, or local loss. Yet reward systems based on such simplified assumptions may not be applicable in any real-world scenario as the dominant strategy for an individual-rational agent is dishonest behavior (e.g., reporting the best possible outcome without costly model training). Recent technologies such as TEE and ZK-SNARKs are promising for trusted computation on mobile, edge, and IoT devices [111]. However, how to leverage them to achieve honest reports without incurring additional costs (e.g., computational costs) is an open question.

Secondly, relative contribution measurement based on the client's reputation or majority voting is an interesting research avenue, promising to relax heavy verification and control mechanics for high-reputation clients. However, how to quantify the reputation fairly and robustly remains open research. Similarly, the majority voting methods may not reflect actual contribution due to its nature.

Thirdly, absolute or direct contribution measurement refers to assessing each client's model update on a public dataset. However, this approach (*i*) requires a trusted central authority performing tests, and (*ii*) limits scalability due to the computational overhead. For instance, the Shapley value is a common method for measuring an agent's contribution, but still comes at the cost of heavy computational overhead even when optimized (e.g., [112], [113]).

Lately, correlation-based reward mechanisms, such as Correlated Agreement (CA) [114], [115] and peer-truth serum [22], have been proposed as promising approaches for contribution measurements for FL. Without having access to the ground truth, the reward is calculated based on the correlation of the reported signals of peers. This implicit approach does not require an explicit contribution measurement and therefore avoids computational overhead.

In addition, when theoretically developing and analyzing incentive mechanisms, as performed by 12 out of 40 papers, more sophisticated assumptions concerning (*i*) information availability (*ii*) uniformity in utility functions or (*iii*) individual rationality should be made to guarantee the robustness of the mechanisms in a real-world scenario. Specifically, as clients are humans, they may not follow their optimal strategies derived from the analysis. For instance, not all clients would take the cost of energy consumption into account when determining their strategies. We suggest taking humans' behavioral bias (e.g., prospect theory [116], [117]) as well as non-quantifiable measures (e.g., the utility of privacy) into the theoretical analysis of incentive mechanisms.

## 6  CONCLUSION AND OUTLOOK

FL is a promising new AI paradigm focused on confidential and parallel model training on the edge. To apply FL beyond small groups of entrusted entities, a decentralization of power, as well as compensation for participating clients, has to be incorporated into the FLF. This work traversed and analyzed 12 leading scientific databases for incentivized and decentralized FLFs based on the PRISMA methodology, ensuring transparency and reproducibility. We found 422 papers and studied 40 works in-depth after three filtering rounds. To ensure correctness, the results were verified by the respective authors. We overcame the challenge of heterogeneity of FLFs in terms of use cases, applied focus, design choice, and thoroughness by offering a comprehensive and holistic comparison framework. By exposing the limitations of existing FLFs and providing directions for future research, this work aims to enhance the proliferation of incentivized and decentralized FL in practice.

## REFERENCES

[1] J. Clement, "Google, Amazon, Facebook, Apple, and Microsoft (GAFAM) - Statistics & facts," *Statista*, 2021.

[2] S. Wheatley, T. Maillart, and D. Sornette, "The extreme risk of personal data breaches and the erosion of privacy," *The European Physical Journal B*, vol. 89, no. 1, p. 7, 2016.

[3] T. Mirrlees, "Getting at GAFAM's power: A structural and relational framework," 02 2021.

[4] C. Santesteban and S. Longpre, "How big data confers market power to big tech: Leveraging the perspective of data science," *The Antitrust Bulletin*, vol. 65, no. 3, pp. 459–485, 2020.

[5] 2018 reform of EU data protection rules. European Commission. [Online]. Available: https://ec.europa.eu/commission/sites/beta-political/files/data-protection-factsheet-changes_en.pdf

[6] H. Brendan McMahan, E. Moore, D. Ramage, S. Hampson, and B. Agüera y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proc. AISTATS*, 2017.

[7] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet *et al.*, "Advances and open problems in federated learning," *arXiv:1912.04977*, 2019.

[8] R. Zeng, C. Zeng, X. Wang, B. Li, and X. Chu, "A comprehensive survey of incentive mechanism for federated learning," *arXiv:2106.15406*, 2021.

[9] J. Hamer, M. Mohri, and A. T. Suresh, "FedBoost: A communication-efficient algorithm for federated learning," in *Proc. of ICML*, vol. 119, 2020, pp. 3973–3983.

[10] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," in *Proc. of Machine Learning and Systems*, vol. 2, 2020, pp. 429–450.

[11] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, "Tackling the objective inconsistency problem in heterogeneous federated optimization," in *Proc. NeurIPS*, vol. 33, 2020, pp. 7611–7623.

[12] F. Sattler, S. Wiedemann, K.-R. Müller, and W. Samek, "Robust and communication-efficient federated learning from non-iid data," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 9, pp. 772–785, 2020.

[13] D. Rothchild, A. Panda, E. Ullah, N. Ivkin, I. Stoica, V. Braverman, J. Gonzalez, and R. Arora, "FetchSGD: Communication-efficient federated learning with sketching," *arXiv:2007.07682*, 2020.

[14] S. Nakamoto, "Bitcoin: A peer-to-peer electronic cash system," *Consulted*, vol. 1, p. 2012, 2008.

[15] S. S. Shetty, C. A. Kamhoua, and L. L. Njilla, "Distributed consensus protocols and algorithms," in *Blockchain for Distributed Systems Security*, 2019, pp. 25–50.

[16] G. Wood *et al.*, "Ethereum: A secure decentralised generalised transaction ledger," *Ethereum Project Yellow Paper*, vol. 151, no. 2014, pp. 1–32, 2014.

[17] "Ethereum networking layer," https://ethereum.org/en/developers/docs/networking-layer/, accessed: 2022-10-18.

[18] "Ethereum proof-of-stake," https://ethereum.org/en/developers/docs/consensus-mechanisms/pos/, accessed: 2022-10-18.

[19] J. Weng, J. Weng, J. Zhang, M. Li, Y. Zhang, and W. Luo, "DeepChain: Auditable and Privacy-Preserving Deep Learning with Blockchain-Based Incentive," *IEEE Transactions on Dependable and Secure Computing*, vol. 18, pp. 2438–2455, 2021.

[20] Y. Li, X. Tao, X. Zhang, J. Liu, and J. Xu, "Privacy-preserved federated learning for autonomous driving," *IEEE Transactions on Intelligent Transportation Systems*, pp. 1–12, 2021.

[21] S. Jiang and J. Wu, "A reward response game in the blockchain-powered federated learning system," *International Journal of Parallel, Emergent and Distributed Systems*, vol. 37, no. 1, pp. 68–90, 2022.

[22] L. Witt, U. Zafar, K. Shen, F. Sattler, D. Li, and W. Samek, "Reward-based 1-bit compressed federated distillation on blockchain," *arXiv:2106.14265*, 2021.

[23] K. Toyoda and A. N. Zhang, "Mechanism design for an incentive-aware blockchain-enabled federated learning platform," in *Proc. of International Conference on Big Data*. IEEE, 2019, pp. 395–403.

[24] K. Toyoda, J. Zhao, A. N. S. Zhang, and P. T. Mathiopoulos, "Blockchain-enabled federated learning with mechanism design," *IEEE Access*, vol. 8, pp. 219 744–219 756, 2020.

[25] Y. Liu, Z. Ai, S. Sun, S. Zhang, Z. Liu, and H. Yu, *FedCoin: A Peer-to-Peer Payment System for Federated Learning*. Springer, 2020, pp. 125–138.

[26] "Matic whitepaper," https://github.com/maticnetwork/whitepaper, accessed: 2022-10-18.

[27] "Binance chain whitepaper," https://github.com/bnb-chain/whitepaper/blob/master/WHITEPAPER.md, accessed: 2022-10-18.

[28] "Avalanche whitepaper," https://assets.website-files.com/5d80307810123f5ffbb34d6e/6008d7bbf8b10d1eb01e7e16_Avalanche%20Platform%20Whitepaper.pdf, accessed: 2022-10-18.

[29] D. G. Wood, "Polkadot: Vision for a heterogeneous multi-chain framework," https://polkadot.network/PolkaDotPaper.pdf.

[30] E. Androulaki, A. Barger, V. Bortnikov, C. Cachin *et al.*, "Hyperledger Fabric: A distributed operating system for permissioned blockchains," in *Proc. of EuroSys Conference*, 2018.

[31] N. Nisan *et al.*, "Introduction to mechanism design (for computer scientists)," *Algorithmic Game Theory*, vol. 9, pp. 209–242, 2007.

[32] S. Chakrabarti, T. Knauth, D. Kuvaiskii, M. Steiner, and M. Vij, "Trusted execution environment with Intel SGX," in *Proc. of Responsible Genomic Data Sharing*, 2020, pp. 161–190.

[33] X. Tu, K. Zhu, N. C. Luong, D. Niyato, Y. Zhang, and J. Li, "Incentive mechanisms for federated learning: From economic and game theoretic perspective," *IEEE Transactions on Cognitive Communications and Networking*, vol. 8, no. 3, pp. 1566–1593, 2022.

[34] H. v. Stackelberg *et al.*, *Theory of the market economy*. Oxford University Press, 1952.

[35] L. U. Khan, S. R. Pandey, N. H. Tran, W. Saad, Z. Han, M. N. H. Nguyen, and C. S. Hong, "Federated learning for edge networks: Resource optimization and incentive mechanism," 2019. [Online]. Available: https://arxiv.org/abs/1911.05642

[36] Y. Zhan, P. Li, Z. Qu, D. Zeng, and S. Guo, "A learning-based incentive mechanism for federated learning," *IEEE Internet of Things Journal*, vol. 7, no. 7, pp. 6360–6368, 2020.

[37] M. Vojnović, "Contest Theory," *Communications of the ACM*, vol. 60, no. 5, pp. 70–80, Apr. 2017.

[38] P. Bolton and M. Dewatripont, *Contract theory*. MIT press, 2004.

[39] V. Krishna, *Auction theory*. Academic press, 2009.

[40] G. Tullock, "Efficient rent seeking," in *Efficient rent-seeking*. Springer, 2001, pp. 3–16.

[41] H. Yu, Z. Liu, Y. Liu, T. Chen, M. Cong, X. Weng, D. Niyato, and Q. Yang, "A Fairness-aware Incentive Scheme for Federated Learning," in *Proc. of the AAAI/ACM Conference on AI, Ethics, and Society (AIES)*, Feb. 2020, pp. 393–399.

[42] Y. Zhan, J. Zhang, Z. Hong, L. Wu, P. Li, and S. Guo, "A survey of incentive mechanism design for federated learning," *IEEE Transactions on Emerging Topics in Computing*, pp. 1–1, 2021.

[43] D. Hou, J. Zhang, K. L. Man, J. Ma, and Z. Peng, "A systematic literature review of blockchain-based federated learning: Architectures, applications and issues," in *Proc. of ICTC*, 2021, pp. 302–307.

[44] A. Ali, I. Ilahi, A. Qayyum, I. Mohammed, A. Al-Fuqaha, and J. Qadir, "Incentive-driven federated learning and associated security challenges: A systematic review," *TechRxiv*, 2021.

[45] D. C. Nguyen, M. Ding, Q.-V. Pham, P. N. Pathirana, L. B. Le, A. Seneviratne, J. Li, D. Niyato, and H. V. Poor, "Federated learning meets blockchain in edge computing: Opportunities and challenges," *IEEE Internet of Things Journal*, vol. 8, no. 16, pp. 12 806–12 825, 2021.

[46] Z. Wang and Q. Hu, "Blockchain-based federated learning: A comprehensive survey," 2021. [Online]. Available: https://arxiv.org/abs/2110.02182

[47] D. Moher, L. Shamseer, M. Clarke, D. Ghersi, A. Liberati, M. Petticrew, P. Shekelle, and L. A. Stewart, "Preferred reporting items for systematic review and meta-analysis protocols (PRISMA-P) 2015 statement," *Systematic Reviews*, vol. 4, no. 1, p. 1, 2015.

[48] C. Okoli and K. Schabram, "A guide to conducting a systematic literature review of information systems research," *SSRN Electronic Journal*, vol. 10, 05 2010.

[49] B. Kitchenham, "Procedures for performing systematic reviews," *Keele, UK, Keele University*, vol. 33, no. 2004, pp. 1–26, 2004.

[50] S. Kumar, S. Dutta, S. Chatturvedi, and M. Bhatia, "Strategies for enhancing training and privacy in blockchain enabled federated learning," in *Proc. of BigMM*, 2020, pp. 333–340.

[51] W. Zhang, Q. Lu, Q. Yu, Z. Li, Y. Liu, S. K. Lo, S. Chen, X. Xu, and L. Zhu, "Blockchain-based federated learning for device failure detection in industrial IoT," *IEEE Internet of Things Journal*, vol. 8, no. 7, pp. 5926–5937, 2021.

[52] Y. Zhao, J. Zhao, L. Jiang, R. Tan, D. Niyato, Z. Li, L. Lyu, and Y. Liu, "Privacy-preserving blockchain-based federated learning for IoT devices," *IEEE IoT Journal*, vol. 8, no. 3, pp. 1817–1829, 2021.

[53] S. Rahmadika and K.-H. Rhee, "Reliable collaborative learning with commensurate incentive schemes," in *Proc. of IEEE International Conference on Blockchain*. IEEE, 2020, pp. 496–502.

[54] Q. Zhang, Q. Ding, J. Zhu, and D. Li, "Blockchain empowered reliable federated learning by worker selection: A trustworthy reputation evaluation method," in *Proc. of WCNCW*, 2021, pp. 1–6.

[55] Z. Zhang, D. Dong, Y. Ma, Y. Ying, D. Jiang, K. Chen, L. Shou, and G. Chen, "Refiner: A reliable incentive-driven federated learning system powered by blockchain," in *VLDB Endowment*, 2021, vol. 14, no. 12, p. 2659–2662.

[56] S. Rahmadika and K.-H. Rhee, "Unlinkable collaborative learning transactions: Privacy-awareness in decentralized approaches," *IEEE Access*, vol. 9, pp. 65 293–65 307, 2021.

[57] M. H. ur Rehman, K. Salah, E. Damiani, and D. Svetinovic, "Towards blockchain-based reputation-aware federated learning," in *Proc. of INFOCOMW*, 2020, pp. 183–188.

[58] S. Rahmadika, M. Firdaus, S. Jang, and K.-H. Rhee, "Blockchain-enabled 5G edge networks and beyond: An intelligent cross-silo federated learning approach," *Security and Communication Networks*, vol. 2021, 2021.

[59] 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.

[60] J. Kang, Z. Xiong, D. Niyato, S. Xie, and J. Zhang, "Incentive mechanism for reliable federated learning: A joint optimization approach to combining reputation and contract theory," *IEEE Internet of Things Journal*, vol. 6, no. 6, pp. 10 700–10 714, 2019.

[61] X. Bao, C. Su, Y. Xiong, W. Huang, and Y. Hu, "FLChain: A blockchain for auditable federated learning with trust and incentive," in *Proc. of BIGCOM*, Aug. 2019, pp. 151–159.

[62] S. Ma, Y. Cao, and L. Xiong, "Transparent contribution evaluation for secure federated learning on blockchain," in *Proc. of ICDEW*, 2021, pp. 88–91.

[63] L. Feng, Z. Yang, S. Guo, X. Qiu, W. Li, and P. Yu, "Two-layered blockchain architecture for federated learning over mobile edge network," *IEEE Network*, pp. 1–14, 2021.

[64] Kansra, Bhrigu and Diddee, Harshita and Sheikh, Tariq Hussain and Khanna, Ashish and Gupta, Deepak and Rodrigues, Joel J. P. C., "BlockFITS: A federated data augmentation modelling for blockchain-based IoVT systems," in *ICICC*, 2022, pp. 253–262.

[65] Mugunthan, Vaikkunth and Rahman, Ravi and Kagal, Lalana, "BlockFLow: Decentralized, privacy-preserving, and accountable federated machine learning," in *BLOCKCHAIN*, 2022, pp. 233–242.

[66] Fadaeddini, Amin and Majidi, Babak and Eshghi, Mohammad, "Privacy preserved decentralized deep learning: A blockchain based solution for secure ai-driven enterprise," in *Proc. of High-Performance Computing and Big Data Analysis*, 2019, pp. 32–40.

[67] C. He, B. Xiao, X. Chen, Q. Xu, and J. Lin, "Federated learning intellectual capital platform," *Personal and Ubiquitous Computing*, 2021.

[68] G. Qu, H. Wu, and N. Cui, "Joint blockchain and federated learning-based offloading in harsh edge computing environments," in *Proc. of the International Workshop on Big Data in Emergent Distributed Environments*, 2021.

[69] H. B. Desai, M. S. Ozdayi, and M. Kantarcioglu, "BlockFLA: Accountable federated learning via hybrid blockchain architecture," in *Proc. of the Eleventh ACM Conference on Data and Application Security and Privacy*, 2021, p. 101–112.

[70] X. Zhu and H. Li, "Privacy-preserving decentralized federated deep learning," in *Proc. of ACM Turing Award Celebration Conference China*, 2021, p. 33–38.

[71] Z. Li, J. Liu, J. Hao, H. Wang, and M. Xian, "CrowdSFL: A secure crowd computing framework based on blockchain and federated learning," *Electronics*, vol. 9, no. 5, 2020.

[72] S. Rathore, Y. Pan, and J. H. Park, "BlockDeepNet: A blockchain-based secure deep learning for IoT network," *Sustainability*, vol. 11, no. 14, 2019.

[73] I. Martinez, S. Francis, and A. S. Hafid, "Record and reward federated learning contributions with blockchain," in *Proc. of International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery (CyberC)*. IEEE, Oct. 2019, pp. 50–57.

[74] Y. Zou, F. Shen, F. Yan, J. Lin, and Y. Qiu, "Reputation-based regional federated learning for knowledge trading in blockchain-enhanced IoV," in *Proc. of IEEE WCNC*, Mar. 2021, pp. 1–6.

[75] Q. Hu, Z. Wang, M. Xu, and X. Cheng, "Blockchain and federated edge learning for privacy-preserving mobile crowdsensing," *IEEE Internet of Things Journal*, pp. 1–1, 2021.

[76] Y. Wang, Z. Su, N. Zhang, and A. Benslimane, "Learning in the air: Secure federated learning for UAV-assisted crowdsensing," *IEEE Trans. Network Sci.Eng.*, vol. 8, no. 2, pp. 1055–1069, Apr. 2021.

[77] M. H. ur Rehman, K. Salah, E. Damiani, and D. Svetinovic, "Towards blockchain-based reputation-aware federated learning," in *Proc. of IEEE INFOCOMW*, Jul. 2020, pp. 183–188.

[78] X. Qu, S. Wang, Q. Hu, and X. Cheng, "Proof of federated learning: A novel energy-recycling consensus algorithm," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 8, pp. 2074–2085, Aug. 2021.

[79] R. Zhang, M. Song, T. Li, Z. Yu, Y. Dai, X. Liu, and G. Wang, "Democratic learning: hardware/software co-design for lightweight blockchain-secured on-device machine learning," *International Journal of High Performance Systems Architecture*, vol. 118, p. 102205, Sep. 2021.

[80] L. Gao, L. Li, Y. Chen, W. Zheng, C. Xu, and M. Xu, "FIFL: A fair incentive mechanism for federated learning," in *Proc. of ICPP*, 2021, pp. 1–10.

[81] S. Xuan, M. Jin, X. Li, Z. Yao, W. Yang, and D. Man, "DAM-SE: A blockchain-based optimized solution for the counterattacks in the internet of federated learning systems," *Security and Communication Networks*, vol. 2021, p. 9965157, Jul. 2021.

[82] Y. Liu, Y. Qu, C. Xu, Z. Hao, and B. Gu, "Blockchain-enabled asynchronous federated learning in edge computing," *Sensors*, vol. 21, no. 10, May 2021.

[83] Y. Liu, X. Yuan, Z. Xiong, J. Kang, X. Wang, and D. Niyato, "Federated learning for 6g communications: Challenges, methods, and future directions," *China Communications*, vol. 17, no. 9, pp. 105–118, sep 2020. [Online]. Available: https://doi.org/10.23919%2Fjcc.2020.09.009

[84] R. Brown, J. Carlyle, I. Grigg, and M. Hearn, "Corda: An introduction," 09 2016.

[85] J. Benet, "IPFS - content addressed, versioned, p2p file system," 2014.

[86] M. Shayan, C. Fung, I. Beschastnikh, and C. J. Yoon, "Biscotti: A ledger for private and secure peer-to-peer machine learning," *arXiv:1811.09904*, 2018.

[87] C. Fung, C. J. M. Yoon, and I. Beschastnikh, "Mitigating Sybils in Federated Learning Poisoning," *arXiv:1808.04866*, Aug. 2018.

[88] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein *et al.*, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.

[89] G. Radanovic, B. Faltings, and R. Jurca, "Incentives for effort in crowdsourcing using the peer truth serum," *ACM Transactions on Intelligent Systems and Technology*, vol. 7, no. 4, pp. 48:1–48:28, 2016.

[90] L. S. Shapley, "A value for N-person games," *Edited by Emil Artin and Marston Morse*, p. 343, 1953.

[91] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for privacy-preserving machine learning," in *Proc. of ACM SIGSAC Conference on CCS*, 2017, p. 1175–1191.

[92] L. Deng, "The MNIST database of handwritten digit images for machine learning research [best of the web]," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.

[93] A. Krizhevsky, V. Nair, and G. Hinton, "The CIFAR-10 dataset," 2014.

[94] M. Everingham, L. van Gool, C. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (VOC) challenge," *IJCV*, vol. 88, no. 2, pp. 303–338, 2010.

[95] M. S. Jere, T. Farnan, and F. Koushanfar, "A taxonomy of attacks on federated learning," *IEEE Security & Privacy*, vol. 19, no. 02, pp. 20–28, mar 2021.

[96] Z. Wang, Q. Hu, and Z. Xiong, "Resource optimization for blockchain-based federated learning in mobile edge computing," 2022. [Online]. Available: https://arxiv.org/abs/2206.02243

[97] A. M. Pinto, "An introduction to the use of zk-SNARKs in blockchains," in *Mathematical research for blockchain economy*. Springer, 2020, pp. 233–249.

[98] H. Yang, J. Zhao, Z. Xiong, K.-Y. Lam, S. Sun, and L. Xiao, "Privacy-preserving federated learning for uav-enabled networks: Learning-based joint scheduling and resource management," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 10, pp. 3144–3159, 2021.

[99] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.

[100] L. U. Khan, S. R. Pandey, N. H. Tran, W. Saad, Z. Han, M. N. Nguyen, and C. S. Hong, "Federated learning for edge networks: Resource optimization and incentive mechanism," *IEEE Communications Magazine*, vol. 58, no. 10, pp. 88–93, 2020.

[101] L. U. Khan, W. Saad, Z. Han, E. Hossain, and C. S. Hong, "Federated learning for internet of things: Recent advances, taxonomy, and open challenges," *IEEE Communications Surveys & Tutorials*, 2021.

[102] K. Guo, S. Han, S. Yao, Y. Wang, Y. Xie, and H. Yang, "Software-hardware codesign for efficient neural network acceleration," *IEEE Micro*, vol. 37, no. 2, pp. 18–25, 2017.

[103] Z. Wang, B. Che, L. Guo, Y. Du, Y. Chen, J. Zhao, and W. He, "PipeFL: Hardware/software co-design of an FPGA accelerator for federated learning," *IEEE Access*, 2022.

[104] S. Caldas, S. M. K. Duddu, P. Wu, T. Li, J. Konečný, H. Brendan McMahan, V. Smith, and A. Talwalkar, "LEAF: A benchmark for federated settings," *arXiv:1812.01097*, Dec. 2018.

[105] W. Wang, "Implementation of federated learning on Raspberry Pi boards: Implementation of federated learning on Raspberry Pi boards with Paillier encryption," 2021.

[106] Lin, He, Zeng, Wang, Huang, and others, "Fednlp: A research platform for federated learning in natural language processing," *arXiv preprint arXiv*, 2021.

[107] M. Servetnyk, C. C. Fung, and Z. Han, "Unsupervised Federated Learning for Unbalanced Data," in *Proc. of Global Communications Conference (GLOBECOM)*. IEEE, Dec. 2020, pp. 1–6.

[108] F. Sattler, A. Marban, R. Rischke, and W. Samek, "Cfd: Communication-efficient federated distillation via soft-label quantization and delta coding," *IEEE Transactions on Network Science and Engineering*, vol. 9, no. 4, pp. 2025–2038, 2022.

[109] T. F. Schranghamer, A. Oberoi, and S. Das, "Graphene memristive synapses for high precision neuromorphic computing," *Nature communications*, vol. 11, no. 1, pp. 1–11, 2020.

[110] H. Yang, K.-Y. Lam, L. Xiao, Z. Xiong, H. Hu, D. Niyato, and H. Vincent Poor, "Lead federated neuromorphic learning for wireless edge artificial intelligence," *Nature communications*, vol. 13, no. 1, pp. 1–12, 2022.

[111] D. Oliveira, T. Gomes, and S. Pinto, "uTango: an open-source TEE for IoT devices," *IEEE Access*, vol. 10, pp. 23 913–23 930, 2022.

[112] R. Jia, D. Dao, B. Wang, F. A. Hubis, N. Hynes, N. M. Gurel, B. Li, C. Zhang, D. Song, and C. Spanos, "Towards efficient data valuation based on the shapley value," *arXiv:1902.10275*, 2020.

[113] T. Wang, J. Rausch, C. Zhang, R. Jia, and D. Song, "A principled approach to data valuation for federated learning," *arXiv:2009.06192*, 2020.

[114] Y. Liu and J. Wei, "Incentives for federated learning: a hypothesis elicitation approach," *arXiv:2007.10596*, 2020.

[115] H. Lv, Z. Zheng, T. Luo, F. Wu, S. Tang, L. Hua, R. Jia, and C. Lv, "Data-free evaluation of user contributions in federated learning," in *19th Symposium on WiOpt*, 2021, pp. 1–8.

[116] D. Kahneman and A. Tversky, "Prospect theory: An analysis of decision under risk," *Handbook of the fundamentals of financial decision making: Part I*, vol. 47, no. 2, pp. 263–292, 1979.

[117] A. Tversky and D. Kahneman, "Advances in prospect theory: Cumulative representation of uncertainty," *Journal of Risk and Uncertainty*, vol. 5, no. 4, pp. 297–323, 1992.

**Wojciech Samek** (M'13) is a professor at the Department of Electrical Engineering and Computer Science at the Technical University of Berlin and is jointly heading the Department of Artificial Intelligence at Fraunhofer Heinrich Hertz Institute (HHI), Berlin, Germany. He received the Master's degree in computer science from the Humboldt University of Berlin, Germany, in 2010, and a Ph.D. degree from the Technical University of Berlin in 2014. He is Associate Faculty with the Berlin Institute for the Foundation of Learning and Data (BIFOLD) and the ELLIS Unit Berlin. He has co-authored more than 150 peer-reviewed publications, several of which were listed by Thomson Reuters as highly cited papers. He is Senior Editor at IEEE TNNLS, and serves on the editorial boards of Pattern Recognition. Furthermore, he is an elected member of the IEEE MLSP Technical Committee and a recipient of multiple best paper awards, including the 2020 Pattern Recognition Best Paper Award and the 2022 Digital Signal Processing Best Paper Prize. His research interest includes deep learning, explainable and thrustworthy AI, and federated learning.

**Leon Witt** is a Ph.D. student at the Department of Computer Science and Technology at Tsinghua University in Beijing. He obtained a Master's degree in Mechanical Engineering and Business Administration from RWTH Aachen, Germany with exchange semesters in Zurich and Los Angeles. He obtained a second Master's degree in Industrial Engineering at Tsinghua University in 2017. His research interests lie at the intersection of federated artificial intelligence, blockchain, and mechanism design.

**Mathis Heyer** is a Master's degree student in Chemical Engineering and Industrial Engineering at RWTH Aachen University, Germany, and Tsinghua University, China. He obtained his Bachelor's degree in Mechanical Engineering from RWTH Aachen University in 2021. As a visiting student, he spent the academic year 2019/2020 at Carnegie Mellon University, USA. His current research interests lie in the applications of artificial intelligence in fields such as chemical engineering and industrial engineering.

**Dan Li** is currently a full professor at the Department of Computer Science and Technology at Tsinghua University-authored the NASP (Network Architecture, System and Protocols) research group, which is part of the networking research lab. He joined the faculty of Tsinghua University in March 2010, after two years working in the Wireless & Networking Group of Microsoft Research Asia as an associate researcher. His main research direction includes trustworthy internet, data center networks and data-driven networking.

**Kentaroh Toyoda** was born in Tokyo, Japan in 1988. He received B.E., M.E., and Ph.D. (Engineering) degrees in the Department of Information and Computer Science, Keio University, Yokohama, Japan, in 2011, 2013, and 2016, respectively. He was an assistant professor at Keio University from Apr. 2016 to Mar. 2019 and is currently a scientist at A*STAR, Singapore. His research interests include blockchain, mechanism design, security and privacy, and data analysis.