

Coalition Formation Game in the Cross-Silo Federated Learning System

Suhan Jiang and Jie Wu

Department of Computer and Information Sciences, Temple University
{tug67249, jiewu}@temple.edu

Abstract—In cross-silo federated learning (FL), organizations cooperatively train a global model with their local data. The organizations, however, own different datasets and may be heterogeneous in terms of their expectation on the precision of the global model. Meanwhile, the cost of secure global model aggregation, including computation and communication, is proportional to the *square* of the number of organizations in the FL system. In this paper, we consider all organizations in the FL system as a grand coalition. We introduce a novel concept from coalition game theory which allows the dynamic formation of coalitions among organizations. A simple and distributed merge and split algorithm for coalition formation is constructed. The aim is to find an ultimate coalition structure that allows cooperating organizations to maximize their utilities in consideration of the coalition formation cost. Through this novel game theoretical framework, the FL system is able to self-organize and form a structured network composed of disjoint stable coalitions. To fairly distribute cost in each formed coalition, a cost sharing mechanism is proposed to align members' individual utility with their coalition's utility. In FL systems, training data has a significant impact on model performances, *i.e.*, it should lead to a more precise global model if organizations with greater data complementarity are grouped. Numerical evaluations are presented to verify the proposed models.

Index Terms—Coalition game, cost sharing, cross-silo federated learning, data quality, horizontal training.

I. INTRODUCTION

Introduced in 2017 [1], federated learning [2] has enabled multiple entities (clients) to collaborate in training a shared model, under the coordination of a central server. Each client uses local data samples without actual exchanging or transferring, and therefore protect data privacy and security. Currently, federated approaches have moved into the mainstream with primary research on extremely large scale settings, composed of millions of mobile and edge devices. Such a *cross-device* FL setting consists of an organization as a model requester and a set of mobile/edge devices as model trainers. The organization acts as the central server to orchestrate the training process and the devices are the clients and perform local training. All clients have no right to make use of the global model since the organization is the *only owner*. Cross-device FL usually involves a huge quantity of clients, each owning a small amount of data.

In recent years, interest in applying FL to a so-called *cross-silo* setting has greatly increased. In this paradigm, there are a small number of relatively reliable clients, each of which represents a larger data store - this setting is more representative of individual companies or organizations (*e.g.*

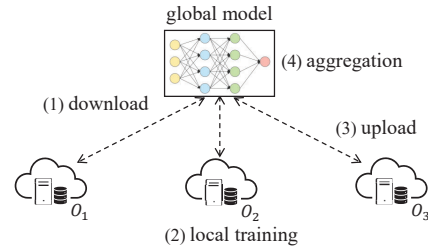


Fig. 1: Cross-silo federated learning system: (1) server sends current global model to all organizations; (2) each organization trains its model using the local data; (3) all organizations upload their updated models to server; (4) server aggregates all local models into a new global model.

financial or medical) with reliable communications and abundant computing resources [3]. Fig. 1 shows a typical cross-silo federated learning system, consisting of a third party entity as the central server and a set of organizations as clients. In this system, the server is responsible for the coordination of training whereas the organizations perform local training. The server first distributes a global model to the organizations, then organizations train the model on locally available data. All updated models are then sent back to the server, where they are averaged to produce a new global model. This new model now acts as the primary model and is again distributed to the organizations. This process is repeated forever or until the global model achieves a satisfactory result from the organization side. Usually, the aggregated global model becomes marginally better than it already was. **All organizations are the co-owners** and can make use of the global model.

In this work, we focus on cross-silo federated learning systems. As we mentioned before, besides local training, federated learning involves interaction between the central server and each client, which is by no means cheap. Costs of model uploading and downloading are inevitable; meanwhile, overhead on the global model aggregation cannot be ignored, especially when secure aggregation [4] is applied. In cross-silo FL system, although some operations, *i.e.*, communication and aggregation, are performed by the central server, the induced costs would be borne by all the participating organizations. To some extent, all of these organizations have formed a coalition, where they collaboratively train a machine learning model with better accuracy, compared to individual training, and share extra costs caused by their cooperation. However, it is still a question whether the grand coalition is stable enough. Is there any chance that some organizations would rather leave the grand coalition and get benefited by forming

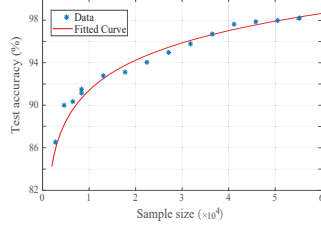


Fig. 2: Relation between the accuracy of the global model and the trained data size: $f(x) = 3.98 \log(9.68 \times 10^5 x - 3.69 \times 10^8)$.

a new group? The answer is affirmative if their utility, *i.e.*, the difference between the model accuracy and the corresponding cost, can be increased in this small coalition. Usually, the model accuracy is positively related to the size of overall trained data, meaning that, the more organizations a coalition contains, the better its global model should be. However, existing works also confirm that the model accuracy and the data size show a concave down increasing trend, indicating the principle of diminishing marginal return. The costs caused by cooperation also increase as the coalition becomes larger, which at least increase linearly or even in a power growth pattern (if applying specific secure aggregation mechanisms). Therefore, when the increase in model accuracy cannot offset the extra cooperation cost due to the large size of the grand coalition, a small coalition would be preferred.

We exploit coalition game theory to analyze the complex interactions among participating organizations. Unlike canonical coalition games which focus on how to stabilize the grand coalition through different reward sharing mechanisms, our work considers a coalition formation game, discussing how to form an appropriate coalition structure, *i.e.*, how organizations are grouped into disjoint stable coalitions, where each small coalition's utility is better than or at least equal to that of the grand coalition. We start from a simple setting where the expected accuracy of a global model is only related to the size of the trained data. We prove the instability of the grand coalition and hence devising a simple and distributed merge-and-split algorithm for forming disjoint coalitions. The coalition-level cost must be shared in a fair manner to promote a stable and long-term cooperation among members. As organizations are grouped based on coalition-level utility, there is a gap between individual utility and coalition utility, leading to free-riding members who want to acquire the global model without contributing local training. To avoid such a situation, we also design a fair cost sharing mechanism to align members' individual utility with the coalition utility. The major contributions of this paper are as follows:

- We propose a coalition formation game to solve a multi-organization grouping problem in a cross-silo federated learning system.
- We show that our proposed game is not superadditive, thereby the grand coalition is seldom the optimal structure.
- We devise a merge-and-split algorithm to form a structured network composed of disjoint stable coalitions.
- We design a cost sharing mechanism with desirable proper-

ties, *e.g.* group-strategyproofness and sharing incentive, to fairly split the cost in each coalition.

- We perform numerical evaluation based on real-world data and the results are consistent with all the theoretical results.

II. PRELIMINARY AND SYSTEM MODEL

A. Cross-Silo Federated Learning

In traditional cross-silo FL, all participating organizations, assuming N in total, aim to cooperatively build a global machine learning model under the orchestration of a central server. First, the server sends the current global model w_G^t to all organizations, where t denotes the current round index. Based on the global model w_G^t , each organization o_i uses its local data to update the local model parameters w_i^t . The goal of organization o_i in round t is to find optimal parameters w_i^t that minimize the loss function $loss(w_i^t)$, *i.e.*, $w_i^{t*} = \arg \min loss(w_i^t)$. Then, each organization uploads its updated local model parameters to the server. Finally, the server facilitates the computation of the parameter aggregation and obtains a new global model w_G^{t+1} . (Several aggregation mechanisms have been proposed for FL [1, 5–7], and here, we focus on Federated-Averaging (FedAvg) [6], in which, the central server updates the global model by summing the weighted models.) We consider that these four steps form a global update round. In a global round, each organization experiences many local training iterations, depending on its training data size. Since the final global model is obtained through many training rounds, here, we only consider one round, where all participating organizations want to improve the updated global model as much as possible. Since we assume that the global model is obtained in a centralized way, the operation costs from the server side, including communication and computation incurred by aggregation, will be distributed among the organizations.

According to the existing works, the accuracy of a machine learning model mainly depends on the training data size. The relation between them can be captured by a concavely increasing function (an example is given in Fig. 2), indicating a decreasing marginal gain. Starting from a simple setting, we assume that all training data in each organization has the same quality and is independently and identically distributed (IID). Based on this assumption, the more data trained by the organizations, the better global model they will obtain at the end of a round. Meanwhile, organizations training more local data make more contributions to the global model. Thus, the server side operating cost should be fairly distributed among organizations based on their individual contributions to the global model. That is, more contribution means less payment.

B. Secure Aggregation for Federated Learning

The security guarantee offered by FL is that sharing updates does not leak any information about the actual training instances used by the clients. Unfortunately, it has been shown that an adversary can invert an individual model update of a target client in order to leak a large amount of information about its local data [8–10]. A central server is the most

Mode	Computation	Communication
Server	$\mathcal{O}(mn^2)$	$\mathcal{O}(mn + n^2)$
Client	$\mathcal{O}(mn + n^2)$	$\mathcal{O}(m + n)$

TABLE I: Computation and communication per round of secure aggregation.

Symbol	Description
O / S	a set of organizations / a coalition(subset) of O
\mathcal{W}	coalition structure, where $\mathcal{W} = \{S_1, \dots, S_k\}$
$N / S $	number of organizations in O / S
o_i	the i -th organization
d_i	data size of o_i
a_i	contribution of o_i in the corresponding coalition
$l(S)$	function measuring the satisfaction level of S
$c(S) / q(S)$	function measuring the cost / data quality of S

TABLE II: Summary of Notations.

vulnerable link, since all local models are directly sent to the server for aggregation. This also means that, individual clients' updates are inspectable by the server. There exist some works [11–13] aiming to improve data privacy in FedAvg. They typically prevent access to the local updates using secret sharing, encryption, and/or reduce information leakage by applying noise to achieve differential privacy. In this paper, we assume that secure aggregation [12] is used to protect client-side privacy. Secure aggregation is a secure multi-party computation protocol that uses encryption to allow a set of clients to compute the sum of their private inputs securely, *i.e.*, only the resulting sum is revealed to the server. Such a protocol requires at least 4 communication rounds between each client and the server in each iteration, which causes extra overhead for the server as well as clients. In practice, this limits the maximum size of a secure aggregation to hundreds of clients. Computational and communication complexities of secure aggregation are listed in Table I, where n is the total number of local models (*i.e.*, number of clients) and m is the length of model updates.

C. A Cross-Silo Federated Learning System

This paper focuses on a cross-silo federated learning system. The model consists of several organizations, aiming to cooperate on model training with their local data. The whole system is in a universal mobile network with wireless communication infrastructures. We consider a quasi-static state where no organizations are joining or leaving. Corresponding notations are shown in Table II.

We consider a scenario with a set O of N organizations, indexed by o_i . Organizations have their own local data while seeking to form cooperative groups, *i.e.*, coalitions, for a better aggregated model. There exists a central orchestration server, which organizes the training for a coalition S , and organizations in S share a common global model. To protect each organization's privacy on the global level, the coalition server applies the secure aggregation protocol. Since there is no limitation on the number of formed coalitions nor a restriction on the size of each coalition, there exist lots of different cooperation methods among these N organizations. Considering an example of 6 organizations, two possible cooperation methods are provided, *i.e.*, method 1: organizations



(a) Method 1: all organizations cooperate as a grand coalition. (b) Method 2: organizations split into two small coalitions.

Fig. 3: A cross-silo federated learning system under two cooperation methods.

form a grand coalition to train a global model (Fig. 3(a)), and method 2: organizations split into 2 small coalitions, each owning a specific global model (Fig. 3(b)).

Organizations have different satisfaction levels when facing different global models. We assume that each organization in the proposed system applies an identical standard to reflect its satisfaction. This standard is defined as the estimated accuracy of the new global model, *i.e.*, a concavely increasing function over the quality of the data trained by all participating organizations. Thus, we use a log function to characterize the relationship between the model accuracy and the training data. Thus, the satisfaction level on the model generated by a coalition S can be expressed as

$$l(S) = \theta \log(1 + \lambda \cdot q(S)), \quad (1)$$

where $q(S)$ is a function used to measure the data quality combined from the coalition S . Based on our previous assumption, $q(S)$ is the total amount of training data in S , *i.e.*, $q(S) = \sum_{o_i \in S} d_i$, where d_i is the size of o_i 's dataset. In reality, there should be more measurements on the quality of combined data. More discussions on explanations of $q(S)$ will be given in Section ??.

Forming a coalition S also incurs cost $c(S)$. In each global secure aggregation, $c(S)$ consists of two parts, *i.e.*, the cost $c_{srv}(S)$ caused by the coalition server and the cost $c_{org}(S)$ from all coalition members. According to Table I, for the coalition server, its costs on both computation and communication grow quadratically with the number of coalition members, denoted as $|S|$. In terms of the organization side, $c_{org}(S)$ is composed of a quadratic of the computation costs and a linear scaling of the communication costs, w.r.t. the number of coalition members. Thus, we simplify the expression of $c(S)$ in Eq. (2):

$$c(S) = c_{srv}(S) + c_{org}(S) = \alpha|S|^2 + \beta|S|. \quad (2)$$

When comparing these two methods in Fig. 3, we could say that, method 1 should provide a better global model for these organizations since more data gathered by the grand coalition, while incurring higher costs as both communication and computation of secure aggregation grows quadratically with the number of clients. Instead, method 2 shows another possibility where each coalition owns a relatively weaker global model, but their summed cost is less than that of forming a grand coalition, meaning less cost borne by each organization. Obviously, since there are many cooperation methods, it is a non-trivial problem to decide which method would be adopted by all the organizations, or how these organizations would reach consensus on a certain method.

It is clear that as the number of organizations per coalition increases, the global model tends to be more accurate while the aggregation cost will increase. This is a crucial trade-off in cross-silo FL that can have a major impact on the collaboration strategies of each organization. Our objective is to derive distributed strategies allowing the organizations to collaborate while accounting for this trade-off.

III. ORGANIZATION COOPERATION: A COALITION FORMATION GAME

To find a suitable partition that satisfies all the organizations, it is natural to consider a centralized approach. In this section, we first formulate this organization cooperation problem from a centralized view, and we show the reason why such a formulation is not suitable in our case. Then, we seek for a distributed solution. Game theory provides a natural paradigm to model the interactions among the organizations in this FL system. Thus, we model the organization cooperation problem as a coalition formation game. Then we prove and discuss its key properties.

A. Centralized Approach

A centralized approach can be used in order to find the optimal coalition structure, that allows the organizations to maximize their benefits from the cross-silo FL. For instance, we seek a centralized solution that maximizes the average model satisfaction level obtained by each organization subject to a cost budget constraint per organization. In a centralized approach, we assume the existence of a centralized entity in the system that is able to gather information on the organizations such as their individual data size d_i or budget b_i . In brief, the centralized entity must be able to know all the required parameters for computing Eq. (1) and Eq. (2) in order to find the optimal structure. Note that, for any organization $o_i \in S$, its individual model satisfaction level is the global model satisfaction level of coalition S , and its budget is limited by the minimal budget in coalition S .

Denoting \mathcal{B} as the set of all partitions of O , the centralized approach seeks to solve the following optimization problem:

Problem 1 ($\text{OP}_{\text{CENTRAL}}$).

$$\text{maximize} \quad \frac{\sum_{S \in \mathcal{P}} |S| \cdot l(S)}{N}, \quad (3a)$$

$$\text{subject to} \quad c(S) \leq |S| \cdot b_S, \quad \forall S \in \mathcal{P}, \quad (3b)$$

where \mathcal{P} is a partition belonging to \mathcal{B} and $b_S = \min_{o_i \in S} b_i$.

Clearly, the centralized optimization problem seeks to find the optimal partition $\mathcal{P}^* \in \mathcal{B}$ that maximizes the average model satisfaction level per organization subject to a budget constraint per coalition. However, it is shown in [14] that finding the optimal coalition structure for solving an optimization problem such as in Problem 1 is an NP-complete problem. This is mainly due to the fact that the number of possible coalition structures (partitions), given by the Bell number, grows exponentially with N , i.e., the number of organization [14]. Moreover, the complexity increases further

due to the fact that the expressions of $l(S)$ and $c(S)$ given by Eq. (1) and Eq. (2) depend on the optimization parameter \mathcal{P} . For this purpose, deriving a distributed solution with a low complexity is desirable. The above formulated centralized approach will be used as a benchmark for the distributed solution in the simulations, for some reasonably small cross-silo FL systems.

B. Game Formulation and Properties

For the purpose of deriving a distributed algorithm that can maximize the model satisfaction level per organization, we refer to cooperative game theory [15] which provides a set of analytical tools suitable for such algorithms. For instance, the proposed organization cooperation problem can be modeled as a (O, u) coalition game [15] where O is the set of players (the organizations) and u is the utility function or value of a coalition. The value $u(S)$ of a coalition $S \subseteq O$ must capture the trade-off between the model accuracy and the aggregation cost. For this purpose, $u(S)$ must be an increasing function of $l(S)$ and a decreasing function of $c(S)$, within coalition S . A suitable utility function is given as below:

$$u(S) = l(S) - c(S) \\ = \theta \log(1 + \lambda \cdot q(S)) - (\alpha |S|^2 + \beta |S|). \quad (4)$$

Traditionally, coalition game based problems seek to characterize the properties and stability of the grand coalition of all players since it is generally assumed that the grand coalition maximizes the utilities of the players. In our case, although forming the grand coalition improves the global model accuracy for the organizations; the cost in terms of aggregation limits this gain. Therefore, for the proposed (O, u) coalition game, we will prove that the grand coalition cannot form due to cost.

Definition 1. A coalition game (O, u) with a transferable utility is said to be superadditive if for any two disjoint coalitions $S_i, S_j \subset O$, $u(S_i \cup S_j) \geq u(S_i) + u(S_j)$.

Theorem 1. The proposed organization cooperation game (O, u) with cost is, in general, non-superadditive.

Proof. Consider two disjoint coalitions $S_i, S_j \subset O$ in the system, their individual utility can be expressed as:

$$u(S_i) = \theta \log(1 + \lambda q_i) - (\alpha |S_i|^2 + \beta |S_i|)$$

$$u(S_j) = \theta \log(1 + \lambda q_j) - (\alpha |S_j|^2 + \beta |S_j|)$$

where $q_i = \sum_{o_k \in S_i} d_k$ and $q_j = \sum_{o_k \in S_j} d_k$.

If coalitions S_i and S_j union, their training data size should be combined as $\sum_{o_k \in S_i \cup S_j} d_k$, i.e., $q_{i \cup j} = (q_i + q_j)$. Thus, the unified utility turns into:

$$u(S_i \cup S_j) = \theta \log(1 + \lambda q_{i \cup j}) - (\alpha |S_i \cup S_j|^2 + \beta |S_i \cup S_j|) \\ \text{The difference between the unified utility and the sum of individual utilities of } S_i \text{ and } S_j \text{ can be found as below:} \\ u(S_i \cup S_j) - [u(S_i) + u(S_j)] \\ = \theta \log(1 + \lambda q_{i \cup j}) - \theta \log[(1 + \lambda q_i)(1 + \lambda q_j)] - 2\alpha |S_i| |S_j| \\ = \theta \log \frac{1 + \lambda q_{i \cup j}}{1 + \lambda q_i + \lambda^2 q_i q_j} - 2\alpha |S_i| |S_j|. \quad (5)$$

size $ S $	number	sum of utilities
0	1	0
1	N	$\sum_{i=1}^N \delta_1^i \theta \log(1 + \lambda d_i)$
2	C_N^2	$\sum_{i=1}^{C_N^2} \delta_2^i [\theta \log(1 + \lambda D_2^i) - (2^2 \alpha + 2\beta)]$
\dots	\dots	\dots
N	1	$\delta_N \theta \log(1 + \lambda \sum_{i=1}^N d_i)$

TABLE III: Sum of utilities under different sizes of coalitions, where $\delta_{|S|}^i$ is the weight for the i -th coalition of size $|S|$.

Since $1 + \lambda q_{i \cup j} < 1 + \lambda q_{i \cup j} + \lambda^2 q_i q_j$, the result of Eq. (5) < 0 holds. Therefore, $u(S_i \cup S_j) < u(S_i) + u(S_j)$; hence the game is not superadditive. \square

Suppose that all organizations cooperate as a grand coalition, then the corresponding utility is denoted as $u(O)$. Denote $\mathbf{x} = (x_1, \dots, x_N)$ as a payoff vector, where x_i is the payoff that organization o_i receives in the grand coalition.

Definition 2. A payoff vector $\mathbf{x} = (x_1, \dots, x_N)$ is said to be group rational or efficient if $\sum_{i=1}^N x_i = u(O)$. A payoff vector \mathbf{x} is said to be individually rational if each organization can obtain the benefit no less than acting alone, i.e., $x_i \geq u(\{o_i\})$, $\forall i$. An imputation is a payoff vector that is group rational and individually rational.

Definition 3. An imputation \mathbf{x} is said to be unstable through a coalition S if $u(S) > \sum_{o_i \in S} x_i$, i.e., the organizations have incentive to form coalition S and reject the proposed \mathbf{x} . The set C of stable imputations is called the core, i.e.,

$$C = \left\{ \mathbf{x} : \sum x_i = u(O) \text{ and } \sum_{o_i \in S} x_i \geq u(S), \forall S \subset O \right\}.$$

A non-empty core means that the organizations have an incentive to form the grand coalition. Let 2^N be the collection of all coalitions, a weighting scheme assigns to every conceivable coalition S a weight $\delta(S)$ between 0 and 1.

Definition 4. A weighting scheme $\delta(\cdot)$ is balanced if it has the property that for every organization o_i , $\sum_{S \in 2^N} \delta(S) = 1$.

Lemma 1. A game (O, u) has a non-empty core if and only if for every balanced weighting scheme $\delta(\cdot)$, $u(N) \geq \sum_{S \in 2^N} \delta(S) u(S)$.

Theorem 2. In general, the core of the proposed (O, u) coalition game is empty.

Proof. Based on Lemma 1, if there exists a balanced weighting scheme $\delta(\cdot)$ that leads to $u(N) < \sum_{S \in 2^N} \delta(S) u(S)$, then we can prove that the proposed (O, u) coalition game has an empty core.

In Table III, we list the sum of utilities under different sizes of coalitions. Here, we consider a special weight scheme δ , where δ_N is x ($0 < x < 1$), and $\delta_1^i = (1 - x)/N$ for $\forall i \in [1, N]$. In this case, we can obtain the following result:

$$\begin{aligned} u(N) - \sum_{S \in 2^N} \delta(S) u(S) &= (1 - x)u(N) - \frac{1 - x}{N} \sum_{i=1}^N u(o_i) \\ &= \frac{1 - x}{N} \theta \left[N \log(1 + \lambda D_N) - \sum_{i=1}^N \log(1 + \lambda d_i) \right] \end{aligned}$$

$$\begin{aligned} &= \frac{1 - x}{N} \theta \left[\sum_{i=1}^N \log \frac{1 + \lambda D_N}{1 + \lambda d_i} \right] - (1 - x)(\alpha N^2 + \beta N) \\ &\leq \frac{1 - x}{N} \theta \left[N \log \frac{1 + \lambda D_N}{1 + \lambda d_{\min}} \right] - (1 - x)(\alpha N^2 + \beta N) \\ &< (1 - x) \left[\theta \log \frac{D_N}{d_{\min}} - (\alpha N^2 + \beta N) \right]. \end{aligned} \quad (6)$$

The data size of each organization in our proposed system should not differ from each other too much, as all of them have sufficient storage and computation. Therefore, in generally, Eq. (6) < 0 holds. Then, we can say Theorem 2 is proven. \square

As a result of the non-superadditivity of the game and the emptiness of the core, the grand coalition does not form among cooperating organizations. Instead, independent disjoint coalitions will form in the system. Therefore, we seek a novel algorithm for coalition formation that accounts for the properties of the organization cooperation game with cost.

IV. DISTRIBUTED COALITION FORMATION ALGORITHM

In this section, we propose a distributed coalition formation algorithm and we discuss its main properties.

A. Coalition Formation Concepts

1) *Orders:* Various criteria (referred to as orders) can be used as comparison relations between partitions, among which, coalition value orders and individual value orders are the most widely-used. Given two partitions \mathcal{W} and \mathcal{P} over the same organization set, where $\mathcal{W} = \{S_1, \dots, S_l\}$ and $\mathcal{P} = \{S'_1, \dots, S'_k\}$, coalition value orders compare the value of two partitions, such as the utilitarian order, in which $\mathcal{W} \triangleright \mathcal{P}$ implies $\sum_{i=1}^l u(S_i) > \sum_{j=1}^k u(S'_j)$. In contrast, the individual value orders compare the individual payoff of each organization, such as the Pareto order. Assuming that the payoff vector in these two partitions are denoted by \mathbf{x} and \mathbf{x}' , which can be written as $\mathcal{W} \triangleright \mathcal{P} \Leftrightarrow \{x_i \geq x'_i, \forall o_i \in \mathcal{W}, \mathcal{P}\}$ by Pareto order if the partition \mathcal{W} is better than the partition \mathcal{P} . In other words, \mathcal{W} is preferred to \mathcal{P} if at least one organization's utility is increased without decreasing other organizations' utilities. Obviously, our proposed system applies utilitarian order.

2) *Stability Notions:* The result of the proposed algorithm in Algorithm 1 is a partition composed of disjoint independent coalitions of organizations. The stability of this resulting structure can be investigated using the concept of a defection function \mathbb{D} [16].

Definition 5. A defection function \mathbb{D} is a function which associates with each partition $\mathcal{P} = \{S_1, \dots, S_k\}$ (each S_i is a coalition) of the player set O , a group of collections in O . A partition \mathcal{P} of O is \mathbb{D} -stable if no group of organizations is interested in leaving \mathcal{P} when the players who leave can only form the collections allowed by \mathbb{D} .

There exist two important defection functions: $\mathbb{D}_{hp}(\mathcal{P})$ (denoted \mathbb{D}_{hp}) and $\mathbb{D}_c(\mathcal{P})$ (denoted \mathbb{D}_c). \mathbb{D}_{hp} associates each partition \mathcal{P} of O with the group of all partitions of O that the players can form through the merge-and-split operation

applied to \mathcal{P} . This function allows any group of players to leave the partition \mathcal{P} of O through merge-and-split operations to create another partition in O . \mathbb{D}_c associates each partition \mathcal{P} of O with the family of all collections in O . This function allows any group of players to leave the partition \mathcal{P} of O through *any operation* and create an arbitrary collection in O . Two forms of stability stem from these definitions: \mathbb{D}_{hp} stability and a stronger \mathbb{D}_c stability. A partition \mathcal{P} is \mathbb{D}_{hp} -stable, if no players in \mathcal{P} are interested in leaving \mathcal{P} through merge-and-split to form other partitions in O ; while a partition \mathcal{P} is \mathbb{D}_c -stable, if no players in \mathcal{P} are interested in leaving \mathcal{P} through any operation (not necessary merge or split) to form other collections in O . Characterizing any type of \mathbb{D} -stability for a partition depends on various properties of its coalitions. For instance, a partition \mathcal{P} is \mathbb{D}_{hp} -stable if, for the partition \mathcal{P} , no coalition has an incentive to split or merge.

Briefly, a \mathbb{D}_{hp} -stable partition can be thought of as a state of equilibrium where no coalitions have an incentive to pursue coalition formation through merge or split. A stronger form of stability can be sought using strict \mathbb{D}_c -stability. The appeal of a strictly \mathbb{D}_c -stable partition is two fold: (1) it is the unique outcome of any arbitrary iteration of merge and split operations done on any partition of O ; (2) it is a partition that maximizes the social welfare, which is the sum of the utilities of all coalitions in a partition. However, the existence of such a partition is not guaranteed. With regards to \mathbb{D}_c -stability, the work in [16–18] proved that a partition $\mathcal{P} = \{S_1, \dots, S_k\}$ of the whole space O is strictly \mathbb{D}_c -stable only if it can fulfill two necessary and sufficient conditions:

- $\forall z \in [1, k]$ and each pair of disjoint coalitions s_i and s_j such that $s_i \cup s_j \subset S_z$, we have $u(s_i \cup s_j) > u(s_i) + u(s_j)$.
- For the partition $\mathcal{P} = \{S_1, \dots, S_k\}$, a coalition $G \subset O$ formed of players belonging to different $S_i \in \mathcal{P}$ is \mathcal{P} -incompatible, that is, $\forall x \in [1, k]$, we have $G \not\subset S_i$. Strict \mathbb{D}_c -stability requires that for all \mathcal{P} -incompatible coalitions G , $\sum_{j=1}^k u(S_i \cap G) > u(G)$.

Therefore, in the case where a partition \mathcal{P} of O satisfying the above two conditions exists; the proposed algorithm converges to this optimal strictly \mathbb{D}_c -stable partition since it constitutes a unique outcome of any arbitrary iteration of merge and split. However, if no such partition exists, the proposed algorithm yields a final network partition that is \mathbb{D}_{hp} -stable.

B. Coalition Formation Algorithm

To ensure autonomous coalition formation, we propose a distributed algorithm based on two simple rules [16], denoted as merge and split, that allow to modify a partition \mathcal{P} of the organizations set O .

Definition 6. Merge Rule - Merge any set of coalitions $\{S_1, \dots, S_k\}$ where $\sum_{j=1}^k u(S_j) < u(\cup_{j=1}^k S_j)$ so that $\{S_1, \dots, S_k\} \rightarrow \cup_{j=1}^k S_j$.

Definition 7. Split Rule - Split any set of coalitions $\cup_{j=1}^k S_j$ where $\sum_{j=1}^k u(S_j) > u(\cup_{j=1}^k S_j)$ so that $\cup_{j=1}^k S_j \rightarrow \{S_1, \dots, S_k\}$.

Algorithm 1 Adaptive Coalition Formation: merge-and-split

Initial: The coalition structure of the network is $\mathcal{P} = \{S_1, \dots, S_N\}$, where $S_i = \{o_i\}$, i.e., all organizations are non-cooperative in the beginning.

Output: an updated coalition structure $\mathcal{P} = \{S_1, \dots, S_k\}$

```

1: repeat
2:   for  $S_i \in \mathcal{W}$  do
3:     Randomly connect to another coalition  $S_j$ 
4:     Perform Merge Rule
5:     Perform Split Rule
6: until merge-and-split terminates
7: Return updated  $\mathcal{P}$ 

```

Based on the merge-and-split rules and the utilitarian order we discussed above, the organizations will form a new coalition if split-and-merge operations improve the utility; otherwise, they keep the original coalition. That is, a group of coalitions decides to merge if it is able to improve its total utility through the merge; while a coalition splits into smaller coalitions if is able to improve the total utility.

A coalition formation algorithm based on merge and split can be formulated for our proposed FL system. Each stage of our coalition formation algorithm will run in two consecutive phases, as shown in Algorithm 1: adaptive coalition formation, and then FL training. During the coalition formation phase, the organizations form coalitions through an iteration of arbitrary merge and split rules repeated until termination. Following the self organization of the system into coalitions, cooperation takes place with each coalition training its own global model. Subsequently, the training phase may occur several times prior to the repetition of the coalition formation phase. It is proven in that any iteration of successive arbitrary merge and split operations terminates.

Theorem 3. Every partition resulting from our proposed coalition formation algorithm is \mathbb{D}_{hp} -stable.

Proof. A partition \mathcal{P} resulting from the proposed merge and split algorithm can no longer be subject to any additional merge or split operation as successive iteration of these operations terminate [18]. Therefore, the organizations in the final partition \mathcal{P} cannot leave this partition through merge and split and the partition \mathcal{P} is immediately \mathbb{D}_{hp} -stable. \square

V. FAIR COST SHARING

The proposed algorithm will yield several coalitions. Moreover, each coalition has its own cost, mainly from secure aggregation. It is necessary to find a way to distribute the cost among all members in the same coalition. Since secure aggregation involves the central server and all members in the coalition, and consumes resources of communication and computation, the simplest sharing method is to divide the cost equally among members. In other words, for organization o_i that belongs to coalition S , its individual utility is:

$$u_i = l(S) - c(S)/|S| - C_i^{lt} = l(S) - (\alpha|S| + \beta), \quad (7)$$

where C_i^{lt} is o_i 's local training cost when it belongs to S . As we mentioned before, a major purpose of our proposed coalition formation algorithm is to seek *coalition-wide* high efficiency in terms of the model accuracy as well as cost. As self-interested and autonomous entities, organizations may behave strategically to maximize their own utility, thereby harming the efficiency. Such an equal division indicates a possibility of organizations being free-riders after joining the coalition, since local training is cost-consuming. For cross-silo FL, the computational and communication resources are non-excludable in the sense that even if an organization does not perform any local training, other organizations cannot exclude that organization from using the trained global model. Thus, we want our cost sharing mechanism to be incentive compatible, *i.e.*, it is in an organization's best interest to perform its local training rather than being free-riders. Also, it should provide an incentive for organizations to participate in the coalition without coercion, *i.e.*, it is fair and maintains the stability of a given coalition formation result.

A. Proportional Fairness

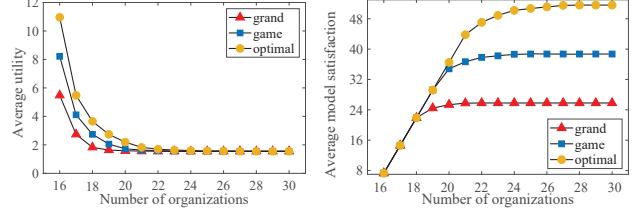
It is reasonable to distribute the coalition-wide cost based on each member's contribution to the coalition model and resource consumption caused by secure aggregation. In fact, the resource consumption of each member performing secure aggregation can be viewed as identical. Similarly, the central server's cost also accumulates equally from communicating and computing with each member. All members' consumption can be treated as equal. Here we focus on individual contribution, denoted a_i . The corresponding organization-side cost will be shared using the function: $c(S)a_i / \sum_{o_j \in S} a_j$. When the coalition cost is distributed based on individual contribution, which reflects the efforts of local training, each organization's individual utility is aligned with its coalition's utility. Such alignment motivates organizations to do the local training as they promised before coalition formation starts, which can effectively avoid free riders.

There exist two common ways to measure an organization's individual contribution. One is using the size of its trained data, and the other is using the organization's local model accuracy. In terms of size-based measurement, a_i can be expressed using $a_i = d_i$. If using accuracy-based measurement, a_i can be expressed using $a_i = \theta \log(1 + \lambda d_i)$. Another measurement, which is widely used in canonical coalition games to maintain the stability of an existing coalition, is Shapely Value. It is a unique mapping ϕ from the coalition utility to individual contribution. For a formed coalition S , ϕ_i , shown in Eq. (8), is the payoff given to organization o_i by the Shapley value ϕ .

$$\phi_i = \sum_{s \subseteq S - \{i\}} \frac{|s|! (|S| - |s| - 1)!}{|S|!} [v(s \cup \{i\}) - v(s)] \quad (8)$$

B. Theoretical Analysis

We present theoretical analysis to demonstrate that our cost sharing mechanism achieves desirable properties. For group-strategyproofness, we should demonstrate that each organization will honestly disclose his real local data size, even if they



(a) Utility of different strategies. (b) Satisfaction of different strategies.

Fig. 4: Impact of organization number in HFL.

are permitted to collude. If a organization's dominant strategy is to truthfully tell the size of its local dataset, then truth-revealing is its dominant strategy.

The cost sharing scheme applied by each coalition S is a function, denoted as ξ , which distributes the total coalition cost $c(S)$ to its members, *i.e.*, ξ takes two arguments, a subset of members Q and an organization o_i , and returns a nonnegative real number satisfying the following: (1) if $o_i \notin Q$ then $\xi(Q, o_i) = 0$, and (2) $\sum_{o_i \in Q} \xi(Q, o_i) = C(Q)$. As is proven in [19], if ξ_j is cross-monotone, then the mechanism specified above is group-strategyproof. Thus, we need to prove ξ_j is cross-monotone. A cost sharing method can be said as cross-monotone if for $Q \subseteq R$, $\xi_j(Q, o_i) \geq \xi_j(R, o_i)$ for every $o_i \in Q$.

Lemma 2. $\xi(S, o_i) = c(S)a_i / \sum_{o_j \in S} a_j$ is cross-monotone.

Proof. Any $o_i \in R \setminus G$ refers to a organization not participating in the FL performed among all organizations in the set of R , thereby they are charged zero cost share. Meanwhile, for all $o_i \in Q \cap R$, $\xi(Q, o_i) = \xi(R, o_i)$. Thus, ξ is a special cross-monotone cost sharing mechanism. \square

Theorem 4. Our cost sharing mechanism satisfies group-strategyproofness and sharing incentive for all organizations.

Proof. The property of group-strategyproofness can be proven using Lemma 2. To show sharing incentive, we should reveal that for any organization, leaving his current assigned coalition would not bring it more benefits. Since our coalition formation algorithm is \mathbb{D}_{hp} -stable, no one has incentive to leave. \square

VI. EVALUATION

Our evaluation will focus on two parts. The first part is to show the advantage of our proposed coalition formation algorithm by comparing with some other coalitional strategies. The second part will analyze the cost sharing mechanism under different definitions of individual contribution. A brief introduction on the experiment settings is given below.

1) *Dataset and Model:* We divide MNIST training samples as organizations' local datasets. We will consider two cases, *i.e.*, equal division and random division. Each organization trains its own multinomial logistic regression model using the Stochastic Gradient Descent (SGD) approach and a coalition-wide global model is obtained by using FedAvg.

2) *Simulation Parameters:* When performing FL, the global model is considered as converged when the loss between two consecutive global rounds is less than 10^{-5} . In a global round, an organization will run local training with 80 epochs with a learning rate of 0.005.

A. Coalition Formation

In this section, we will conduct two experiments under the total number of organizations, *i.e.*, N , changes. In the first experiment, we set $N = 8$, and we will compare our proposed split-and-merge partition strategy with the optimal partition using the centralized approach (Section III.A). And in the second experiment, N varies in the range of $[15, 30]$, and we only compare our proposed split-and-merge partition strategy with the grand-coalition strategy due to the complexity of the centralized approach. When comparing different partition strategies, we will use (1) the utility of all organizations, (2) average model satisfaction level over all formed coalitions, and (3) average model accuracy over all formed coalitions if FL is really performed, to measure their performances.

In the first experiment where $N = 8$ and $(\theta, \lambda, \alpha, \beta)$ is set as $(10, 8 \times 10^{-6}, 0.05, 0.2)$, we equally distribute MNIST training data so that each organization holds a dataset of 7500 samples. The centralized optimal strategy gives a grand coalition structure of $\mathcal{W}_c = \{\{o_1, o_2, o_3, o_4, o_5, o_6, o_7, o_8\}\}$ if there is no limitation on each organization's budget. In this case, the average model satisfaction level is 14.8 and the total cost is 4.8. However, our split-and-merge partition strategy can yield a solution of 2 coalitions as $\mathcal{W}_d = \{\{o_1, o_2, o_3, o_6\}, \{o_4, o_5, o_7, o_8\}\}$ or a solution of 4 coalitions $\mathcal{W}'_d = \{\{o_1, o_2\}, \{o_3, o_6\}, \{o_4, o_8\}, \{o_5, o_7\}\}$. In the 2-coalition structure, each coalition's utility is 10.8, given the model satisfaction level of 12.4 and the coalition-wide cost of 1.6. In the 4-coalition structure, each coalition's utility is 10.8 as well, while the model satisfaction level is 11.2 and the coalition-wide cost of 0.4. Obviously, the centralized solution targets on a high model satisfaction level while the our distributed solution takes the cost as an important factor for the coalition formation. We further conduct FL under these three structures to see the real model accuracy and convergence time. The average model accuracy resulted from \mathcal{W}_c , \mathcal{W}_d and \mathcal{W}'_d is 99.31, 98.73, and 98.14, respectively. Obviously, those results are aligned with our model satisfaction level. We notice that, for a structure with multiple coalitions, the convergence time of each coalition is close, indicating that little concern on the fairness of waiting time among different coalitions.

Next, we show the average utility and the average model satisfaction level in Fig 4. The grand-coalition strategy always leads to a better model satisfaction level while the cost for communication is a big concern with the increasing N . Obviously, our split-and-merge strategy can achieve a reasonable balance between the model accuracy and the total communication cost.

B. Cost Sharing

In Table IV, we show the impact caused by different definitions on the individual contribution. Obviously, the increase of N 's value will lead to higher average cost as the secure aggregation always increases quadratically as the size of a coalitions and more organizations indicate larger sizes of formed coalitions. Even with different cost sharing mechanism, we can still observe the impacts caused by different

Strategy \ N	10	20	30	40
optimal	16	32	48	62
game	23	34.5	48.4	64
grand	27	38	50	68.1

(a) Average cost under size-based policy.

Strategy \ N	10	20	30	40
optimal	15.8	30.2	44.3	61
game	17.8	31.3	46.9	63.1
grand	24.6	35.9	48.8	64.7

(b) Average cost under accuracy-based policy.

Strategy \ N	10	20	30	40
optimal	15.1	28.8	41.5	51.8
game	17.7	28.9	42.7	52.8
grand	23.9	34.5	43.1	53.9

(c) Average cost under SV-based policy.

TABLE IV: Impact of different definitions on the individual contribution.

coalition formation strategies. Note that, there is no specific way to measure which definition is better. All of these three definitions can be applied as long as the organizations in the system reach an agreement.

VII. RELATED WORK

1) *Federated Learning*: As there is more and more attention on privacy, federated learning has become one of the essential concepts in modern machine learning. The salient feature of FL enables its widespread applications in both cross-device and cross-silo settings. In cross-device FL, more clients are enthralled to contribute their resources to improve their user experience. For example, Google applies FL to its products Gboard to improve the performance [20]. Similarly, Apple employs FL to QuickType. Besides that, FL also demonstrates its potential to solve the dilemma problem of *isolated data island* faced by companies/organizations who hesitate to share their vast volume of data samples for business concerns and privacy regulations [21].

2) *Coalition Game Theory*: Unlike noncooperative game theory that studies competitive scenarios, cooperative game theory provides analytical tools to study the behavior of rational players when they cooperate. As a main branch of cooperative games, coalition games describe the formation of cooperating groups of players, that can strengthen the players' positions in a game [15]. According to [22], coalition games can be grouped into three categories: canonical coalition games, coalition formation games, and coalition graph games. In canonical games [23], the grand coalition composed of all users should be an optimal structure, and major topics in this field focus on how to stabilize the grand coalition [24]. Canonical games implicitly assume that forming a coalition is always beneficial, while formation game [25, 26] admits the presence of a cost for forming coalitions, thereby the coalition structure that forms depends on gains and costs from cooperation. In a graph game [27], there exists a graph representing the connectivity of the players among each other, *i.e.*, which player communicates with which one inside each

and every coalition. However, in canonical and formation games, players are assumed to be fully connected.

3) *Fair Cost-sharing Mechanism*: There exist a number of fair cost-sharing mechanisms [28] for coalition formation games, e.g. equal-split, proportional-split, and egalitarian-split solutions. These mechanisms model practical cost-sharing applications with desirable properties, such as the existence of a stable coalition structure with a small strong price-of-anarchy (SPoA) [29] to approximate the social optimum. [30] devises practically cost-sharing mechanisms for decentralized coalition formation that can lead to desirable stable coalition structures. The challenge of our problem is the gap between the coalition utility and its members' individual utility.

VIII. CONCLUSION

In this paper, we propose a coalition formation game to solve a multi-organization grouping problem in a cross-silo federated learning system. We show that our proposed game is not superadditive, thereby the grand coalition is seldom the optimal structure. A simple and distributed merge and split algorithm for coalition formation is constructed. The aim is to find an ultimate coalition structure that allows cooperating organizations to maximize their utilities while also accounting for the cost of coalition formation. To fairly distribute cost in each formed coalition, a cost sharing mechanism is proposed to align members' individual utility with their coalition's utility. The experimental results show that our scheme is efficient in terms of cost reduction for both the group as a whole and individuals.

REFERENCES

- [1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017.
- [2] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, 2020.
- [3] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, "Advances and open problems in federated learning," *arXiv preprint arXiv:1912.04977*, 2019.
- [4] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for federated learning on user-held data," *arXiv preprint arXiv:1611.04482*, 2016.
- [5] L. Muñoz-González, K. T. Co, and E. C. Lupu, "Byzantine-robust federated machine learning through adaptive model averaging," *arXiv preprint arXiv:1909.05125*, 2019.
- [6] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer, "Machine learning with adversaries: Byzantine tolerant gradient descent," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 118–128.
- [7] D. Yin, Y. Chen, R. Kannan, and P. Bartlett, "Byzantine-robust distributed learning: Towards optimal statistical rates," in *International Conference on Machine Learning*. PMLR, 2018.
- [8] B. Hitaj, G. Ateniese, and F. Perez-Cruz, "Deep models under the gan: information leakage from collaborative deep learning," in *Proceedings of the 2017 ACM SIGSAC Conference on Ccs*, 2017.
- [9] M. Nasr, R. Shokri, and A. Houmansadr, "Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning," in *2019 IEEE symposium on SP*. IEEE, 2019.
- [10] L. Melis, C. Song, E. De Cristofaro, and V. Shmatikov, "Exploiting unintended feature leakage in collaborative learning," in *2019 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2019.
- [11] P. Mohassel and P. Rindal, "Aby3: A mixed protocol framework for machine learning," in *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, 2018.
- [12] 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 *proceedings of the 2017 ACM SIGSAC Conference on CCS*, 2017.
- [13] S. Truex, N. Baracaldo, A. Anwar, T. Steinke, H. Ludwig, R. Zhang, and Y. Zhou, "A hybrid approach to privacy-preserving federated learning," in *Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security*, 2019.
- [14] T. Sandholm, K. Larson, M. Andersson, O. Shehory, and F. Tohmé, "Coalition structure generation with worst case guarantees," *Artificial intelligence*, 1999.
- [15] R. B. Myerson, *Game theory: analysis of conflict*. Harvard university press, 1997.
- [16] K. R. Apt and T. Radzik, "Stable partitions in coalitional games," *arXiv preprint cs/0605132*, 2006.
- [17] K. R. Apt and A. Witzel, "A generic approach to coalition formation," *International game theory review*, 2009.
- [18] G. Demange and M. Wooders, *Group formation in economics: networks, clubs, and coalitions*. Cambridge University Press, 2005.
- [19] N. R. Devanur, M. Mihail, and V. V. Vazirani, "Strategyproof cost-sharing mechanisms for set cover and facility location games," *Decision Support Systems*, 2005.
- [20] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, "Federated learning for mobile keyboard prediction," *arXiv preprint arXiv:1811.03604*, 2018.
- [21] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2019.
- [22] W. Saad, Z. Han, M. Debbah, A. Hjørungnes, and T. Basar, "Coalitional game theory for communication networks," *Ieee signal processing magazine*, 2009.
- [23] R. J. La and V. Anantharam, "A game-theoretic look at the gaussian multiaccess channel," *Tech. Rep.*, 2003.
- [24] C. Singh, S. Sarkar, A. Aram, and A. Kumar, "Cooperative profit sharing in coalition-based resource allocation in wireless networks," *IEEE/ACM Transactions on Networking*, 2011.
- [25] R. J. Aumann and J. H. Dreze, "Cooperative games with coalition structures," *International Journal of game theory*, 1974.
- [26] T. Arnold and U. Schwalbe, "Dynamic coalition formation and the core," *Journal of economic behavior & organization*, 2002.
- [27] J. Derks, J. Kuipers, M. Tennekes, and F. Thuijsman, "Local dynamics in network formation," in *Proc. Third World Congress of The Game Theory Society*. Citeseer, 2008.
- [28] Y. Zhou and S. C.-K. Chau, "Multi-user coalition formation for peer-to-peer energy sharing," in *Proceedings of the Eleventh ACM International Conference on Future Energy Systems*, 2020.
- [29] N. Andelman, M. Feldman, and Y. Mansour, "Strong price of anarchy," *Games and Economic Behavior*, 2009.
- [30] S. C.-K. Chau, K. Elbassioni, and Y. Zhou, "Approximately socially-optimal decentralized coalition formation," *arXiv preprint arXiv:2009.08632*, 2020.