

# Share: Stackelberg-Nash based Data Markets

Jinfei Liu, Yuran Bi, Chen Zhao, Junyi Zhao, Kui Ren

Zhejiang University

{jinfeiliu, stellabyr, zhaochen49, junyizhao, kuiren}@zju.edu.cn

Li Xiong

Emory University

lxiong@emory.edu

**Abstract**—With the prevalence of data-driven intelligence, data markets with various data products are gaining considerable interest as a promising paradigm for commoditizing data and facilitating data flow. In this paper, we present Stackelberg-Nash based Data Markets (*Share*) to first realize the incentive data market construction with absolute pricing for a demand-driven market. We propose a three-stage Stackelberg-Nash game to model trading dynamics which not only optimizes the profits of all selfish participants but also adapts to the common *buyer-broker-sellers* market flow and solves the seller selection problem based on sellers' inner competition. We define Stackelberg-Nash Equilibrium and use backward induction to solve the equilibrium. For inner Nash equilibrium, we apply the conventional direct derivation approach and propose a novel mean-field based method along with provable approximation guarantees for complicated cases where direct derivation fails. Experiments on real datasets verify the effectiveness and efficiency of *Share*.

## I. INTRODUCTION

Data products (e.g., query services, aggregate statistics, and machine learning models) have paved the way for a variety of data-driven tasks in diverse industries. High-performance data products require a large amount of high-quality data. While there are a wealth of data generated from different sources, they are highly dispersed, which brings significant challenges to data aggregation. Besides, there is a gap between data supply and demand, and data suppliers or demanders usually lack the necessary resources and techniques to survey the vast data sources and turn data into data products. Thus, despite the increasingly available and enriched data, the wealth of data is far from being fully exploited. As one of the most important topics in Boston Database Meeting 2023 [1], data markets have been demonstrated as a promising paradigm to commoditize data and connect data suppliers and demanders [2], [3].

**Motivations.** A typical data (product) market consists of three parties: buyers, brokers, and sellers [4]–[6]. Buyers propose demands for data products and pay for them; brokers facilitate the transactions between buyers and sellers (and take charge of manufacturing data products from data); sellers offer data with different quality and sell data to brokers in exchange for compensation. We use two motivating examples to further specify our targeted settings.

*Vehicle Example.* An automaker (e.g., Ford Motors Company) wants to get insight into users' purchase preferences of vehicles to decide investments. Ford (**buyer**) turns to McKinsey, a consulting company, and proposes a series of queries  $Q$ , each with distinctive conditions specified and the corresponding purchase intention to be answered, e.g., *whether a female customer in Texas would most likely buy*

*fuel vehicles, pure electric vehicles, or hybrid vehicles.* To answer the queries, McKinsey (**broker**) needs to gather sales data of different vehicle types and produce a data product  $M$ , e.g., *data statistics aggregating total sales of each vehicle type after screening location and gender.* McKinsey buys data from multiple vehicle retailers (**sellers**) who own sales data of various vehicles. Retailers sell data with different qualities responding to different compensations.

*Health Example*<sup>1</sup>. A biopharmaceutical company, e.g., Pfizer (who spends \$12 million to buy health data from a variety of sources including IMS Health as reported by Scientific America [7]), wants to get insight into the effects of their released COVID-19 vaccination for further development. Pfizer (**buyer**) turns to a healthcare consulting firm with a series of queries  $Q$  searching for the arising adverse reactions, e.g., *select the reported nausea within three months after the vaccination in America.* To answer the queries, the consulting firm (**broker**) needs to gather realistic health data to produce a data product  $M$ , e.g., *data aggregation listing the symptomatic description of nausea with the time and location of vaccine inoculation filtered.* The consulting firm collects data from healthcare companies (**sellers**), e.g., the aforementioned IMS Health who owns and sells de-identified prescription data, medical claims, and electronic medical records. While the health data is fiercely anonymized to comply with privacy laws, IMS Health may still suffer from risks of medical disputes and thus enhance privacy preservation which contributes to the provided health data with different quality.

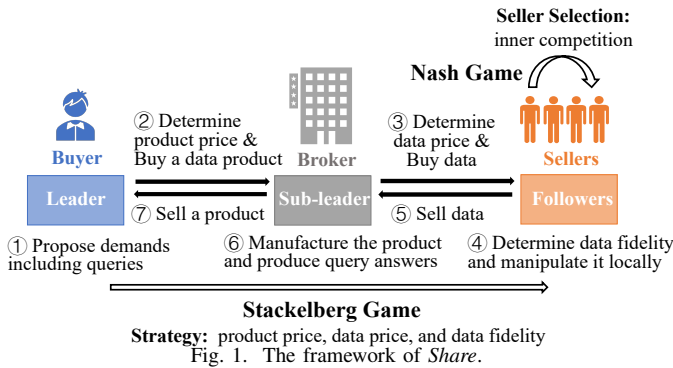
Ford or Pfizer can benefit from the data product (accessing it via querying and answering) and meantime needs to pay for it; McKinsey or the healthcare consulting firm gains by selling the product (answering the queries) after spending resources to buy the data and produce the product; and the vehicle retailers or the healthcare companies sell data for compensations while suffering from costs (majorly the privacy loss incurred from data). They are driven to join the data market by the profit they can earn, leading us to considering a general data market where all three parties are *selfish*, i.e., have their own *revenue* and *cost*, and aim to maximize their own profit (the difference between *revenue* and *cost*). Moreover, how they act in the market affects each other. If Ford sets a low price for the demanded product in pursuit of profit, McKinsey may pay little to buy the data to recover costs. Getting low compensations,

<sup>1</sup>While data sharing in healthcare can bring societal benefits, many issues such as privacy and ethics need to be considered. Here we aim to show the motivation of *Share* supported by real cases, leaving the other issues outside the scope of the paper.

the vehicle retailers offer poor-quality data, inducing a low-performance product and in turn harming the profit of Ford. Therefore, **it's a critical research issue in data management to design an incentive mechanism (especially a data pricing mechanism) for data market which can encourage three selfish yet interdependent parties to participate in data trading and thus invigorate the market.**

While all three parties need to maximize their profit, they take different roles in the market flow. Data sellers such as the vehicle retailers and healthcare companies more likely do not regard data selling as the main business. Rather, data transactions are probably initiated by data buyers (demanders) such as Ford and Pfizer as well as in many real-life scenarios. In such demand-driven scenarios, a single transaction serves for one specific demand (e.g., personalized services) and thus buyers can be considered as orientating the market in turn (coming one at a time) as practiced in work [4]. The data buyer proposes her demand to a data broker in the market, e.g., McKinsey, which then buys data from data sellers. Concerning the limited data owned by one single seller, multiple data sellers are considered, e.g., numerous vehicle retailers in our example, each having a dataset that can together contribute to solving the buyer's demand. Therefore, to facilitate data trading in scenarios like the illustrated examples, **we focus on demand-driven data markets with one data buyer and multiple data sellers.**

Many recent works [4]–[6], [8], [9] have proposed market paradigms emphasizing different aspects of the data market, yet not targeted the demand-driven one with incentives for all three parties. Therefore, it's tempting to ask: **how to build a well-functioning data market with an incentive pricing mechanism, which can satisfy the profit needs of all selfish participants and adapt to the demand-driven scenarios.**



**Challenges and Contributions.** We summarize three challenges ( $C_1, C_2, C_3$ ) faced in constructing the demand-driven data market with incentives for all three parties, and propose a feasible solution, **Stackelberg-Nash based Data Markets (*Share*)** utilizing game theory as in Figure 1. The detailed workflow is presented in Section III-B.

Existing works on data markets vary in design goals and typically address one aspect or one party's need, such as product quality optimization for buyers [10], revenue maximization for sellers [5], social welfare maximization [11], or market protection from strategic participants [12], but fail to realize

profit optimization for all parties. Pricing constitutes a key mechanism in incentive market construction, but the lack of market practices and pricing references makes data pricing far from trivial (similar with petroleum pricing at the early stage), especially an absolute price compared to the relative one which is determined in comparison with other data, e.g., by Shapley value [13] as in [4], [6], [9]. Therefore, the first challenge is ( $C_1$ ): **How to design a pricing mechanism for data markets to realize absolute pricing and to maximize the profits of all three parties.** To solve this challenge, we adopt game theory which can support the multi-objective incentive mechanism design in data markets. The interactions of the three entities are modeled as a game, in which each participant can achieve her profit-maximization goal by making her optimal strategy. Moreover, absolute prices of data are modeled as strategies of participants and directly determined in the game process with the involvement of all the parties, which further encourage their participation.

Though efforts have been made to satisfy buyers' needs (e.g., utility demand and purchase budget in [4], [6]), no existing paradigm can well adapt to the demand-driven data market with specific market flow in order, i.e., first demanded and initiated by the buyer, then translated and transmitted through the broker, and finally received and realized by the sellers. Therefore, the second challenge is ( $C_2$ ): **How to encode the buyer-broker-sellers market flow to cater to the targeted demand-driven data market.** To solve this challenge, we formalize the interactions among three parties as a multi-stage dynamic game and adapt Stackelberg game [14] which can deal with the sequential order of participants, by regarding the buyer as the leader, the broker as the sub-leader, and sellers as followers. As shown in Figure 1, the buyer first announces what data product she demands for and determines the product price based on her profit-maximization goal; the broker then tries to buy data from sellers and decides the data price; each seller then chooses what data quality to provide.

Since there are multiple sellers, it's critical to select the *best* data (with the highest data quality) from the sellers to meet the buyer's product demand and meantime satisfy the broker's resource constraints, which is referred as *seller selection problem*. Many existing works made the buyer [10] or broker [6], [15] responsible for seller selection, which not only requires the buyer/broker's capability of learning the data quality but also limits the sellers' ability of choosing their provided data quality according to the data price. Hence, the third challenge is ( $C_3$ ): **How to model the seller selection problem to select the best set of data for trading.** To solve this challenge, we consider the inner competition among sellers which can make the winners as the selected sellers without the assistance of the buyer or broker. Sellers are allowed to manipulate their provided data quality by Local Differential Privacy (LDP) [16] to compete for the selling quantity of data. We model the inter-seller competition as a Nash game [17] because of its advantage in modeling sellers' equal positions, and find the desired Nash equilibrium by applying direct derivation and proposing a mean-field based

approximation for complex cases when direct derivation fails.

Our goal is not to cover all data markets nor to address all critical issues in real-world data trading, but rather to propose an incentive mechanism for data markets anchored in a demand-driven scenario, which is meaningful in practice but not studied yet. The major contributions are summarized.

- We present *Share*, an incentive data market framework with absolute pricing mechanism based on a three-stage Stackelberg-Nash game, which is the first to satisfy all-party profit maximization in demand-driven scenarios.
- We apply Nash game for seller selection problem, which formulates sellers' inner competition and incorporates seller selection into the game process among three parties.
- We define Stackelberg-Nash Equilibrium in data markets and derive it by backward induction. To solve inner Nash game, we apply direct derivation as well as design a novel mean-field method for complex cases, for which error analysis is presented.
- We conduct experiments on real and synthetic datasets to verify the effectiveness and efficiency of *Share*.

**Organization.** Section II provides the related work. Section III presents our data market framework based on Stackelberg-Nash game. In Section IV, a market instance is constructed, for which approaches to deriving the equilibrium and the trading dynamics are presented. Section V reports the experimental results and findings. Section VI draws a conclusion and discusses future work.

## II. RELATED WORK

We discuss related work on data market and game theory.

### A. Data Market

Data markets trade data in direct or indirect forms (derived data products). Marketplaces where buyers directly purchase raw data have been practiced [18], [19]. [8], [20] proposed query-based data markets which allow buyers to obtain information by querying the database. Recently, model-based data markets [4]–[6], [9] have been proposed where machine learning models trained by data are traded. While research on data market evolves in myriad research directions including data mining [21]–[26], data storage [27], [28], and data security [29], [30], we explore the data market design emphasizing on the incentives by formalizing trading (pricing) mechanisms. Related research problems are reviewed below.

In terms of profit maximization for all parties, few studies can provide a thorough solution. [6] established a model marketplace with the needs of buyers and sellers considered, but assumed that the broker is neutral without her own profit consideration and determines model prices only for single goal optimization, i.e., revenue maximization for sellers. [15], who studied transactions for crowdsensing data, were devoted to multiple goal optimization. Nevertheless, the specific characteristics of crowdsensing data (e.g., sensing time) limit the extension to the general data market. Diversified data products, typical privacy issues, and latent interrelations among participants should be enhanced. In *Share*, by combining multiple

game mechanisms, we formulate for-all profit-maximization data markets with unrestricted data and data products, privacy consideration for data sellers, as well as inner competition modelling for seller selection.

As for data pricing, several surveys [31]–[33] claimed fundamental principles and reviewed the evolution of pricing models. In terms of absolute pricing, [6] provided absolute prices for data models, which, however, highly rely on the survey results and can't be adjusted dynamically. [4] applied Myerson's payment rule to determine absolute model payment but allocated relative compensations to sellers in proportion to their contributions based on Shapley value. While auction [34] can be a promising way for absolute price discovery and has been widely adopted in data pricing [11], [12], rarely can every party be included in the price determination. In *Share*, we propose a feasible absolute pricing mechanism for both data and data product with all three parties involved.

Many works looked at the seller (data) selection problem. One strand of research lied in data acquisition. For example, [10] dug into how a buyer purchases data under a budget to improve machine learning models, yet without touching the market design issues including data pricing, revenue allocation, as well as the strategic actions of sellers and brokers. Within the data market design scope, [6] made brokers choose datasets to maximize Shapley coverage of the trained model. [15] used a combinatorial multi-armed bandit mechanism for brokers to select sellers. However, the selection results directly affect the profits of sellers, and therefore the seller selection problem is closely correlated to the profit maximization problem for sellers and should not be considered separately. In fact, seller selection can be seen as the spontaneous process of the inner competition among sellers, proactively determined by their strategies rather than passively conducted by the buyer or broker. In *Share*, the seller selection problem is formalized as the inter-seller Nash game, which is a part of the incentive mechanism for profit optimization of all participants.

### B. Game Theory

(Non-cooperative) Game theory provides a tool for analyzing the interplay among individuals with conflicting objectives and has been widely used in various situations. Nash [17] accurately described Nash equilibrium as a solution concept for simultaneous-move games. Many researchers used Nash game as a powerful tool to formulate and solve problems with simultaneous interactions [35], [36]. Instead, Stackelberg game [14] features sequential actions, which was first used to formulate the determining process for oligopoly firms producing homogeneous products and has been further applied to many practical situations with hierarchical organizations, e.g., security game [37] and crowdsensing [38].

Since Nash proposed his theory, many researchers have sought algorithms for finding Nash equilibrium. [39] showed complexity results of deriving Nash equilibrium and [40] further studied the complexity of computing a mixed Nash equilibrium. In terms of solving Stackelberg game, backward induction approach, an iterative technique to derive dynamic

game equilibrium, is often used [15], [38], [41]. In fact, deriving Stackelberg equilibrium with complete information can be formulated as a bilevel optimization problem [42].

In this paper, we adopt Stackelberg game for the focused demand-driven data markets because it captures the sequential actions of participants and can thus adapt to the *buyer-broker-sellers* market flow while remaining the profit maximization for all parties. Moreover, we first adopt Nash game for the seller selection problem since Nash game models the simultaneous-move interaction among equals and can be used for the inner competition among data sellers, which can select sellers based on their strategies.

### III. MARKET FRAMEWORK: PARTICIPANTS, MECHANISM, AND EQUILIBRIUM

We first describe the market framework from the perspective of participants in Section III-A. We formulate the market mechanism as a three-stage Stackelberg-Nash game in Section III-B and define the market equilibrium, Stackelberg-Nash Equilibrium in Section III-C. For reference, Table I summarizes the frequently used notations.

TABLE I  
THE SUMMARY OF FREQUENTLY USED NOTATIONS.

	Notation	Definition
Buyer $\mathcal{B}$	$\mathbf{Q}$	demand query to be solved
	$\nu$	demand product performance
	$p^M$	basic price of data product
	$\theta_1, \theta_2$	parameters of concern on each attribute
	$\rho_1, \rho_2$	parameters of sensitivity to each attribute
	$\mathbf{U}(\cdot)$	utility function of the product
	$\mathbf{PB}(\cdot)$	payment function between buyer and broker
Broker $\mathcal{A}$	$\Phi(\cdot)$	profit function of the buyer
	$N$	total data quantity
	$p^D$	basic price of data
	$\sigma_k$	parameters of manufacturing cost
	$\mathbf{C}(\cdot)$	cost function of manufacturing data product
Seller $\mathcal{S}_i$	$\Omega(\cdot)$	profit function of the broker
	$i$	index of seller
	$m$	total number of sellers
	$\tau_i$	data fidelity
	$\epsilon_i$	parameter in local differential privacy
	$\chi_i$	sold data quantity
	$\lambda_i$	parameter of privacy sensitivity
	$\mathbf{L}_i(\cdot)$	privacy loss function
Data	$\mathbf{PS}_i(\cdot)$	payment function between broker and seller
	$\Psi_i(\cdot)$	profit function of the seller
	$D_i$	seller $\mathcal{S}_i$ 's raw dataset
	$D_i^t$	seller $\mathcal{S}_i$ 's provided dataset
	$D^t$	whole dataset for manufacturing
	$q_i^D$	dataset quality provided by seller $\mathcal{S}_i$
	$q^D$	total quality of dataset for manufacturing
	$q^M$	data product quality
	$\omega_i$	weight of seller $\mathcal{S}_i$ 's dataset

#### A. Market Participants

*Share* describes a demand-driven market with a data buyer, a data broker, and data sellers, whose roles are defined below.

- **Buyer.** To fulfill a data-driven task, buyer  $\mathcal{B}$  asks for a data product from broker  $\mathcal{A}$ . Buyer  $\mathcal{B}$  gets access to the product via queries  $\mathbf{Q}$  (either transactional or analytical) and claims her required product performance (the answer

accuracy<sup>2</sup>), notated as  $\nu$ . Note that the trading would fail if the demand of buyer  $\mathcal{B}$  is not satisfied, e.g., the product cannot give answers for the asked 100 queries, or some answer has an accuracy lower than expected. Buyer  $\mathcal{B}$  gains utility from the product (query answers) while giving the payment to broker  $\mathcal{A}$ .

- **Broker.** Broker (Arbiter)  $\mathcal{A}$  wants to make profits by bridging the transactions between buyers and sellers. To answer the queries  $\mathbf{Q}$  with demand  $\nu$  satisfied, broker  $\mathcal{A}$  needs to buy  $N$  data records (limited by her computation resources) from sellers to make the data product (can be in any needed form from statistics aggregating to model training). The product manufacturing incurs certain costs (e.g., computing cost). Then, broker  $\mathcal{A}$  sells the product to buyer  $\mathcal{B}$  in exchange for payment.
- **Sellers.** A large number of sellers  $\{\mathcal{S}_i | i = 1, 2, \dots, m\}$  exist in the market. Each seller  $\mathcal{S}_i$  owns dataset  $D_i$  and wants to sell it for profit. Seller  $\mathcal{S}_i$  applies perturbation to her data utilizing a privacy scheme to manipulate the data quality locally and sells  $\chi_i$  processed data records to broker  $\mathcal{A}$ . The data quantity  $\chi_i$ , with  $\sum_{i=1}^m \chi_i = N$ , is to be decided by the market mechanism, and for any required number  $\chi_i \in \mathbb{N}^+$ ,  $|D_i| \geq \chi_i$ , indicating that each seller has enough data for the trading. Seller  $\mathcal{S}_i$  receives compensation from broker  $\mathcal{A}$  while suffering from the (privacy) cost for the data she sells.

**Information Structure.** Complete information [14] is assumed in Stackelberg game and Nash game. To simplify the model and focus on the formulation of interactions among three profit-seeking parties, we put emphasis on a highly transparent market with public profit functions. For those cases where participants are unwilling to make profit details open to the public due to privacy concerns, we provide an alternative solution by introducing a trusted marketplace which receives all the participants' private profit parameters and takes charge of executing the trading mechanism. In both cases, all the parties are assumed honest to give individual parameters (as illustrated in Section IV-C) under the supervision of market regulators, e.g., enforced by regular spot-check.

**Problem Statement.** Each of the participants has *revenue* (gained utility, received payment, or compensation) and *cost* (payment, manufacturing cost, or privacy loss). All the participants are profit-driven and want to maximize their own profit, i.e., the difference between *revenue* and *cost*. The problem is to find an optimal strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$ , where buyer  $\mathcal{B}$ , broker  $\mathcal{A}$ , and each seller  $\mathcal{S}_i$  determine product price  $p^{M*}$ , data price  $p^{D*}$ , and data fidelity  $\tau_i^*$  respectively such that the profits of all participants are maximized. The profits of the buyer, the broker, and each seller are defined below, including

<sup>2</sup>We use *accuracy* in a broad sense which, independent of our model, can be measured in distinctive ways depending on the query task (e.g., the precision of the returned purchase intention on vehicles, or the completeness of the aggregated nausea cases).

the function templates and desired properties, which can be instantiated in various forms. A market instance with every function instantiated will be constructed in Section IV and feasible approaches are proposed to solve its equilibrium.

*1) Profit Function of Buyer:* When buyer  $\mathcal{B}$  comes to the market and asks for a data product, she cares about her *revenue*, the utility she can get from the product, and her *cost*, the payment she should give to the broker.

*Revenue.* The revenue of buyer  $\mathcal{B}$  is the utility gained from the product. Apparently, the performance of product itself (embodied in the accuracy of query answers which is specified in demand  $\nu$ ) affects the utility. Moreover, the quality of data used to make the product contributes to the utility. While the answer accuracy only indicates how the product performs under a certain testing environment (related to specific validation datasets), dataset quality measures how good *raw materials* are, making the judgment of product utility more stable and less sensitive to various application scenarios. The dataset quality is measured as the total quality of datasets contributed by all sellers,  $q^D = \sum_{i=1}^m q_i^D$ , where  $q_i^D$  is the dataset quality provided by seller  $\mathcal{S}_i$ . Besides the intrinsic characteristic of the data (e.g., the number of data features, completeness, and currency), the dataset quality is positively correlated with the data fidelity  $\tau_i$  provided by seller  $\mathcal{S}_i$  as well as the number of data  $\chi_i$  bought from her, notated as  $q_i^D = g(\chi_i, \tau_i)$  which will be instantiated in Section IV-A. Data fidelity  $\tau_i$  is determined by the perturbation added by seller  $\mathcal{S}_i$ , measured as the privacy level of LDP mechanism (see more in Section III-A3) while  $\chi_i$  is determined by sellers' inner competition on  $\tau$  (see more in Section III-B). Combining both dataset quality and product performance, the gained utility is quantified by a utility function  $\mathbf{U}(q^D, \nu)$  following the law of diminishing marginal utility [43] in economics, as instantiated in Section IV-A.

*Cost.* The *cost* of buyer  $\mathcal{B}$  is the payment to broker  $\mathcal{A}$ . Based on the above analysis, we define  $q^M = h(q^D, \nu)$  to objectively represent the quality of the data product which depends on both data quality  $q^D$  and product performance  $\nu$ , and  $h(\cdot)$  will be instantiated in Section IV-A. Also,  $p^M$  is defined as the basic price of  $q^M$  (the product price), and the payment for the product can be formulated as the function  $\mathbf{PB}(p^M, q^M)$  positively correlated to the basic price and the product quality, which will be instantiated in Section IV-A.

*Profit.* The profit  $\Phi(\cdot)$  of buyer  $\mathcal{B}$  is the difference between the quantification of product utility and the payment to the broker.

$$\Phi(p^M, \tau) = \mathbf{U}(q^D, \nu) - \mathbf{PB}(p^M, q^M). \quad (1)$$

*2) Profit Function of Broker:* When broker  $\mathcal{A}$  receives the demand from buyer  $\mathcal{B}$ , she cares about her *revenue*, i.e., the payment from buyer  $\mathcal{B}$ , and her *cost* consisting of the compensations to sellers to buy the data and the manufacturing cost in the process of producing the data product.

*Revenue.* The *revenue* of broker  $\mathcal{A}$  is the payment from buyer  $\mathcal{B}$ , i.e.,  $\mathbf{PB}(p^M, q^M)$  (the *cost* of buyer  $\mathcal{B}$ ).

*Cost.* The *cost* of broker  $\mathcal{A}$  is the sum of 1) the compensations to sellers and 2) the manufacturing cost. Broker  $\mathcal{A}$  needs to pay each seller  $\mathcal{S}_i$  compensation according to her provided data quality, which is formulated as function  $\mathbf{PS}(p^D, q_i^D)$  and instantiated in Section IV-A. Here  $p^D$  describes the basis price of data (the data price) similar to the product price  $p^M$ . Broker  $\mathcal{A}$  also needs to consume some resources to make the product. Different manufacturing consumption would be induced if processing data with different sizes or producing product with different performance. Therefore, cost function  $\mathbf{C}(N, \nu)$  is formulated related to total data size  $N$  and product performance  $\nu$  and will be instantiated in Section IV-A.

*Profit.* The profit  $\Omega(\cdot)$  of broker  $\mathcal{A}$  is defined as the received payment from buyer  $\mathcal{B}$  minus the compensations to sellers and the manufacturing cost as follows.

$$\Omega(p^M, p^D, \tau) = \mathbf{PB}(p^M, q^M) - \sum_{i=1}^m \mathbf{PS}(p^D, q_i^D) - \mathbf{C}(N, \nu). \quad (2)$$

*3) Profit Function of Seller:* When seller  $\mathcal{S}_i$  gets the purchase request for data from broker  $\mathcal{A}$ , she cares about her *revenue*, the compensation from broker  $\mathcal{A}$  and her *cost* coming mostly from her privacy loss.

*Revenue.* The *revenue* of seller  $\mathcal{S}_i$  is the compensation from broker  $\mathcal{A}$ , i.e.,  $\mathbf{PS}(p^D, q_i^D)$  (one part of the *cost* of broker  $\mathcal{A}$ ).

*Cost.* The *cost* of seller  $\mathcal{S}_i$  is mainly the privacy loss incurred based on data fidelity  $\tau_i$  she provides (we ignore other costs of collecting, processing, and packaging data which can be formulated as a constant and easy to cope with). In our settings, data fidelity  $\tau_i$  is determined and manipulated by seller  $\mathcal{S}_i$  through LDP mechanism which provides a well-justified measurement tool to simultaneously capture noise level (fidelity) and privacy level (cost). Hence,  $\tau_i$  is defined as  $f(\epsilon_i)$  where  $\epsilon_i$  represents the privacy level in standard LDP. We conclude the following characteristics  $f(\cdot)$  should satisfy concerning the marginal trend and boundary conditions, and the instantiation will be shown in Section IV-A.

1. The data has fidelity  $\tau_i = 0$  when  $\epsilon_i = 0$  which means very heavy perturbation has been added, causing the nearly random data.
2. Larger  $\epsilon_i$ , higher  $\tau_i$ , since less noise is added to data.
3.  $\tau_i$  increases slower as  $\epsilon_i$  becomes larger because very little noise is being added and further decreasing noise does not make a significant difference to data fidelity anymore. On the other hand, when  $\epsilon_i$  is very small, i.e., with extremely large noise, increasing  $\epsilon_i$  can significantly increase data fidelity. Besides,  $\tau_i$  cannot increase perpetually and should be upper bounded.

Bigger  $\tau_i$  means better fidelity of data and more privacy loss for seller  $\mathcal{S}_i$ . We quantify such loss by function  $\mathbf{L}_i(\cdot)$  which is positively related to  $\tau_i$ . It's intuitive that the cost function should not only increase but also increase faster for higher  $\tau_i$ , which corresponds to the principle of increasing marginal cost [43] in economics. Moreover, the privacy cost would increase as more data is sold (larger  $\chi_i$ ). Specific function  $\mathbf{L}_i(\tau_i)$  will be elaborated in Section IV-A.

*Profit.* The profit  $\Psi_i(\cdot)$  of seller  $\mathcal{S}_i$  is the difference between the compensation and the quantification of privacy loss.

$$\Psi_i(p^D, \tau_i) = \mathbf{PS}(p^D, q_i^D) - \mathbf{L}_i(\tau_i). \quad (3)$$

### B. Market Mechanism

In *Share*, the three entities take strategies in order. We first present the market workflow. Then we specify the strategies of buyer  $\mathcal{B}$ , broker  $\mathcal{A}$ , and each seller  $\mathcal{S}_i$ , respectively. Based on the strategies, the market mechanism is proposed.

**Market Workflow.** The market workflow is shown in Fig. 1. ① Buyer  $\mathcal{B}$  puts forward the demand for a product including the queries and the required performance. ② Buyer  $\mathcal{B}$  determines the product price to buy the data product from broker  $\mathcal{A}$ . ③ Broker  $\mathcal{A}$ , acting as the bridge for the transaction between the buyer and  $m$  sellers, determines the data price to buy the data from sellers. ④ Each seller chooses what data (strictly speaking, data fidelity) to sell, and conducts corresponding privacy perturbation locally. ⑤ Sellers sell the protected datasets to broker  $\mathcal{A}$  in exchange for the compensations. ⑥ Using the dataset bought from sellers, broker  $\mathcal{A}$  manufactures the product. ⑦ Broker  $\mathcal{A}$  sells the product to buyer  $\mathcal{B}$ . After buyer  $\mathcal{B}$  receives the product via query answers and gives payment to broker  $\mathcal{A}$ , the transaction is finished.

**Buyer's Strategy.** Buyer  $\mathcal{B}$  makes her strategy first, which is to determine the product price  $p^M$ , in order to maximize her profit by considering the desired utility of the product and stimulating the responses of the broker and sellers, i.e., what data price and data fidelity broker  $\mathcal{A}$  and sellers would provide according to  $p^M$ .

**Broker's Strategy.** Broker  $\mathcal{A}$  takes her strategy second, which is to determine data price  $p^D$ , in order to maximize her profit given  $p^M$  by stimulating the sellers' responses, i.e., what data fidelity each seller would provide according to  $p^D$ .

**Seller's Strategy.** Sellers make their strategies last. The strategy of each seller  $\mathcal{S}_i$  is to determine data fidelity  $\tau_i$  to maximize her profit by balancing the revenue of selling data and the cost of the privacy loss given the data price  $p^D$ .

Meanwhile, the inner competition among  $m$  sellers should be considered. Given the data price  $p^D$ , if seller  $\mathcal{S}_i$  provides data with higher fidelity  $\tau_i$ , more quantity would likely be sold. If other sellers provide better fidelity, less data quantity of seller  $\mathcal{S}_i$  could be chosen. The data quantity that each seller  $\mathcal{S}_i$  can sell is thus formalized as  $\chi_i(\tau_i, \tau_{-i})$  (see the instance in Section IV-A). Each seller competes for the quantity of data that can be sold by manipulating the data fidelity while balancing the compensation and the privacy cost. We define such inner competition among sellers as a Nash game. Seller  $\mathcal{S}_i$  determines her strategy  $\tau_i$  simultaneously with each other to maximize her own profit which is also affected by other sellers' strategies  $\tau_{-i}$ . Nash equilibrium would be achieved where no seller can increase her profit by unilaterally changing her strategy with all other sellers' strategies fixed. The data quantity  $\chi_i$  sold by each seller  $\mathcal{S}_i$  can be calculated according to the equilibrium state, treated as the seller selection results.

Note that if one participant finds that her maximized profit is below zero, she will quit since she can gain no benefit from participating in the data trading, which guarantees the individual rationality [44] of participants. If it is the buyer or the broker who quits the trading or all sellers simultaneously get negative profits and quit, the current transaction would fail and the new transaction would be initiated. Otherwise, the remaining participants would continue and finish the transaction. Since it's easy to deal with the quit situation, we focus on the more common case and assume that all participants can get non-negative profits in the following discussions.

**Three-Stage Stackelberg-Nash Game.** Strategies of buyer  $\mathcal{B}$ , broker  $\mathcal{A}$ , and sellers  $\mathcal{S}_i$  ( $i = 1, 2, \dots, m$ ) constitute the strategy profile  $\langle p^M, p^D, \tau \rangle$  of data markets. Such a profile determines market trading rules including selling at what price for both data product ( $p^M$ ) and data ( $p^D$ ), what data (data fidelity) to sell ( $\tau$ ), as well as how to select sellers (the calculated  $\chi = (\chi_1, \chi_2, \dots, \chi_m)$  based on  $\tau$ ). The market mechanism is formulated as a three-stage Stackelberg-Nash game, where buyer  $\mathcal{B}$  is the leader, broker  $\mathcal{A}$  is the sub-leader, and  $m$  sellers act as the followers. Each of them tries to maximize her own profit by determining her optimal strategy variable. The three-stage Stackelberg-Nash game is defined as follows.

*Definition 1 (Three-Stage Stackelberg-Nash Game):* The game consists of three stages for buyer, broker, and sellers.

*Stage 1* Buyer  $\mathcal{B}$ :  $p^{M*} = \arg \max_{p^M} \Phi(p^M, \tau(p^D(p^M)))$ .

*Stage 2* Broker  $\mathcal{A}$ :  $p^{D*} = \arg \max_{p^D} \Omega(p^M, p^D, \tau(p^D))$ .

*Stage 3* Seller  $\mathcal{S}_i$ :  $\tau_i^* = \arg \max_{\tau_i} \Psi_i(p^D, \tau), i = 1, 2, \dots, m$ .

The above three-stage Stackelberg-Nash game involves both sequentiality and simultaneity. Sequentiality indicates the order in market flow, i.e., driven by demand, the data trading proceeds with buyer  $\mathcal{B}$  acting first, broker  $\mathcal{A}$  taking her strategy second, and sellers making their strategies last. Simultaneity indicates the equal positions of  $m$  sellers who take strategy simultaneously in their inner Nash game.

### C. Market Equilibrium

In the above game, our objective is to find an optimal strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$ , by which each participant can maximize her own profit. Meanwhile, the optimal solution must satisfy some equilibrium so that no one is willing to adopt other strategies, which indicates market stability, making our design feasible. We define a Stackelberg-Nash Equilibrium (SNE) in data markets.

*Definition 2 (Stackelberg-Nash Equilibrium):* An optimal strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$  constitutes a Stackelberg-Nash Equilibrium (SNE) if and only if the following set of inequalities is satisfied.

$$\Phi(p^{M*}, \tau^*(p^{D*}(p^{M*}))) \geq \Phi(p^M, \tau^*(p^{D*}(p^M))), \quad (4)$$

$$\Omega(p^M, p^{D*}(p^M), \tau^*(p^{D*})) \geq \Omega(p^M, p^D, \tau^*(p^D)), \quad (5)$$

$$\Psi_i(p^D, \tau^*(p^D)) \geq \Psi_i(p^D, \tau_{-i}^*(p^D), \tau_i), i = 1, 2, \dots, m. \quad (6)$$

SNE indicates that each participant takes her optimal strategy which maximizes her own profit in a demand-driven data market with *buyer-broker-sellers* sequence. No one can add her own profit by unilaterally changing her strategy.

#### IV. MARKET CONSTRUCTION: EQUILIBRIUM SOLVING AND TRADING DYNAMICS

In this section, we first instantiate a data market by specifying each function template in Section IV-A, and then derive the market equilibrium by backward induction in Section IV-B. We describe the market dynamics in Section IV-C.

##### A. Market Instance

In terms of profit functions of participants, we claim the basic properties which should be satisfied in Section III, based on which we give instances below following certain practices. Other alternatives can be adopted based on real cases.

1) *Profit instantiation of Buyer B*: In terms of the utility of buyer  $\mathcal{B}$ , it has been analyzed that combining product performance and dataset quality to measure product utility can make the quantification of product utility more comprehensive. Since the dataset quality  $q_i^D$  of each seller is positively correlated with  $\chi_i$  and  $\tau_i$ , we instantiate  $q_i^D = g(\chi_i, \tau_i)$  as  $\chi_i \tau_i$  for simplicity and thus  $q^D = \sum_{i=1}^m \chi_i \tau_i$ . As mentioned before, other inherent factors of data may also contribute to the data quality, which has been studied yet complementary to our work [45], and we focus on the effect sellers can exert on the data instead of the intrinsic characteristics which, on the other hand, can be further formulated into  $q_i^D$  as a constant. Based on the instance of  $q^D$ , we define the utility of a data product as the weighted sum of the utility of the dataset quality and the utility of the product performance, which are further formulated as the logarithmic functions following utility theory [46] in economics.

$$U(q^D, \nu) = \theta_1 \ln(1 + \rho_1 q^D) + \theta_2 \ln(1 + \rho_2 \nu). \quad (7)$$

Here  $\theta_1$  and  $\theta_2$  satisfy  $\theta_1, \theta_2 \in (0, 1), \theta_1 + \theta_2 = 1$ , which measure the relative significance of the two for buyer  $\mathcal{B}$ . In our example, if dataset quality  $q^D$  plays a greater role than product performance  $\nu$  in the decision-making of the automaker, the automaker may set  $\theta_1 = 0.7$  and  $\theta_2 = 0.3$ .  $\rho_1 > 0$  and  $\rho_2 > 0$  refer to buyer  $\mathcal{B}$ 's sensitivity to these two attributes respectively. More sensitive, more utility added when the attribute gets better. For example, if higher dataset quality can bring the automaker much more utility, its  $\rho_1$  would be big, meaning that the automaker is highly sensitive to the quality of production materials.

In terms of the payment of buyer  $\mathcal{B}$ , we instantiate  $q^M = h(q^D, \nu)$  as  $q^D \nu$  since it is positively correlated to  $q^D$  and  $\nu$ , and specify  $\mathbf{PB}(p^M, q^M) = p^M q^M$ , which borrows from common sense that the payment for goods is equal to the unit (basic) price multiplied by the quantity (quality).

2) *Profit instantiation of Broker A*: In terms of the first part of the cost of broker  $\mathcal{A}$ , the payment to each seller  $\mathcal{S}_i$  is similarly instantiated as the product of the basic price  $p^D$  and the dataset quality  $p^D$ , i.e.,  $\mathbf{PS}(p^D, q_i^D) = p^D q_i^D$ .

In terms of the second part of the cost of broker  $\mathcal{A}$ , we adopt a widely used transcendental logarithmic function for the manufacturing cost because of its adaptability to varied economies of scale and manufacturing strategy (e.g., how to allocate computing resources) according to the work [47]. Here  $\sigma_k, k \in \{0, 1, 2, 3, 4, 5\}$  are the parameters of the translog cost function which can be fitted by broker  $\mathcal{A}$  based on the actual manufacturing procedure.

$$C(N, \nu) = \exp \left( \sigma_0 + \sigma_1 \ln(N) + \sigma_2 \ln(\nu) + \frac{1}{2} \sigma_3 \ln^2(N) + \frac{1}{2} \sigma_4 \ln^2(\nu) + \sigma_5 \ln(N) \cdot \ln(\nu) \right). \quad (8)$$

3) *Profit instantiation of Seller  $\mathcal{S}_i$* : In terms of the privacy cost of seller  $\mathcal{S}_i$ , we first instantiate data fidelity  $\tau_i$  and data quantity  $\chi_i$ . We choose an inverse trigonometric function form as  $f(\cdot)$  which satisfies the characteristics stipulated in Section III-A3 and give the following definition of  $\tau_i$ .

$$\tau_i = f(\epsilon_i) = \frac{2}{\pi} \operatorname{arcsec}(w_i \epsilon_i + 1), \quad \epsilon_i \in [0, \infty), \quad (9)$$

which leads to  $\tau_i \in [0, 1)$ . Additionally,  $\tau_i = 1$  when no noise is added. Thus  $\tau_i \in [0, 1]$ .

We then instantiate the quantity  $\chi_i$  of data that can be sold by seller  $\mathcal{S}_i$  as proportional to the data fidelity  $\tau_i$  she provides.

$$\chi_i = N \frac{\omega_i \tau_i}{\sum_{j=1}^m \omega_j \tau_j}, \quad (10)$$

where  $\omega_1, \omega_2, \dots, \omega_m$  refer to the weights of sellers' data, which are maintained by the broker. Such weights reflect the historical performance of each seller's data in past deals (implying the reputation of sellers). The broker would update these weights after each round of transactions. For example, new weights can be updated based on the contributions of sellers to the data product in the current transaction. One of the evaluation methods for the data contribution is by Shapley value [13], which is adopted in this market instance and implemented in our experiments.

Based on the instances of data fidelity  $\tau_i$  and data quantity  $\chi_i$ , we adopt a widely used quadratic function for the cost of seller  $\mathcal{S}_i$ . Here  $\lambda_i > 0$  is seller  $\mathcal{S}_i$ 's privacy sensitivity. In our health example, the privacy loss of IMS Health corresponds to the negative impact of data exposure and  $\mathbf{L}_i(\cdot)$  quantifies the economic estimation of the impact (e.g., legal expenses or amends to patients).

$$\mathbf{L}_i(\tau_i) = \lambda_i (\chi_i \tau_i)^2. \quad (11)$$

##### B. Solving Equilibrium: Backward Induction

To determine the optimal strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$ , we adopt the backward induction approach [48]. We first investigate Stage 3 to solve Nash equilibrium among sellers and derive the expression of each seller's optimal strategy  $\tau_i^*, i = 1, 2, \dots, m$  (Eq. 15) for any given data price  $p^D$  in Section IV-B1. We explore two methods, direct derivation and an approximate method using the mean-field state which can

deal with complicated cases. Next, we consider Stage 2 to determine the expression of the optimal strategy  $p^{D*}$  (Eq. 20) of broker  $\mathcal{A}$  for any given product price  $p^M$  in Section IV-B2. In this process, the expression of  $\tau_i^*, i = 1, 2, \dots, m$  solved from Nash game can be used as sellers' optimal reactions to  $p^D$ . Then, we back to Stage 1 to find the value (rather than the expression) of buyer  $\mathcal{B}$ 's optimal strategy  $p^{M*}$  (Eq. 22) based on the optimal reactions of the broker as well as sellers in Section IV-B3. After that, we can get the value of the optimal strategy  $p^{D*}$  by substituting  $p^{M*}$  into the result (Eq. 20) in Stage 2. Finally, we can compute the value of each seller's optimal strategy  $\tau_i^*$  by substituting  $p^{D*}$  into the result (Eq. 15) in Stage 3. Till now, the complete strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$  has been determined. The detailed deduction is presented as follows.

1) *Expression of  $\tau^*$  in Stage 3:* We present two approaches to derive the expression of  $\tau_i^*$  for sellers, direct derivation and a mean-field based approximation method for large numbers of sellers and complicated profit function forms which can hardly be solved by direct derivation.

**Direct Derivation.** By substituting Eqs. 10, 11 into Eq. 3, we get each seller's profit

$$\begin{aligned} \Psi_i(p^D, \tau_i) &= p^D \chi_i \tau_i - \lambda_i (\chi_i \tau_i)^2 \\ &= p^D \cdot N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} - \lambda_i \left( N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \right)^2. \end{aligned}$$

$\Psi_i$  is correlated to not only seller  $\mathcal{S}_i$ 's strategy  $\tau_i$  but also other sellers' strategies  $\tau_j, j \neq i$  because of the inner competition formulated as Nash game among sellers. As we discussed before, each seller aims to maximize her own profit. Therefore, we derive each of the first-order derivatives for  $m$  sellers' profit functions and let each of them equal to zero, thus getting  $m$  equations. The equation for seller  $\mathcal{S}_i$  is

$$p^D \frac{\partial \left( N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \right)}{\partial \tau_i} - 2\lambda_i \cdot N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \cdot \frac{\partial \left( N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \right)}{\partial \tau_i} = 0. \quad (12)$$

If  $\frac{\partial \left( N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \right)}{\partial \tau_i} = 0$ , it is an all-zero solution, which does not meet our problem situation, so we can directly eliminate  $\frac{\partial \left( N \frac{\omega_i \tau_i^2}{\sum_{j=1}^m \omega_j \tau_j} \right)}{\partial \tau_i}$ , and then get

$$\begin{cases} p^D \sum_{i=1}^m \omega_i \tau_i - 2N \lambda_1 \omega_1 \tau_1^2 = 0 \\ p^D \sum_{i=1}^m \omega_i \tau_i - 2N \lambda_2 \omega_2 \tau_2^2 = 0 \\ \vdots \\ p^D \sum_{i=1}^m \omega_i \tau_i - 2N \lambda_m \omega_m \tau_m^2 = 0, \end{cases} \quad (13)$$

where each  $\mathcal{S}_i$ 's equation not only relates to her own strategy  $\tau_i$  but also contains other sellers' strategies, requiring us to solve  $m$  simultaneous equations together. Finding that

$$2N \lambda_1 \omega_1 \tau_1^2 = 2N \lambda_2 \omega_2 \tau_2^2 = \dots = 2N \lambda_m \omega_m \tau_m^2 = p^D \sum_{i=1}^m \omega_i \tau_i. \quad (14)$$

By adding all  $m$  equations in Eq. 13, we get

$$mp^D \sum_{i=1}^m \omega_i \tau_i - 2N \sum_{i=1}^m \lambda_i \omega_i \tau_i^2 = 0.$$

Using  $\tau_1$  to indicate other  $\tau_i$  ( $i = 2, 3, \dots, m$ ) from Eq. 14,

$$mp^D \tau_1 \sum_{i=1}^m \sqrt{\frac{\lambda_1 \omega_1 \omega_i}{\lambda_i}} - 2N m \lambda_1 \omega_1 \tau_1^2 = 0.$$

Therefore,

$$\tau_1^* = \frac{p^D}{2N \sqrt{\omega_1 \lambda_1}} \sum_{i=1}^m \sqrt{\frac{\omega_i}{\lambda_i}},$$

and using Eq. 14 again, we get all sellers' optimal strategies

$$\tau_i^* = \frac{p^D}{2N \sqrt{\omega_i \lambda_i}} \sum_{j=1}^m \sqrt{\frac{\omega_j}{\lambda_j}}, i = 1, 2, \dots, m. \quad (15)$$

Note that the second-order derivative  $\frac{\partial^2 \Psi_i(p^D, \tau_i)}{\partial \tau_i^2} < 0$ , so these solutions can maximize each seller's profit.

**Mean-field based Approximate Method.** It's theoretically feasible that the optimal  $\tau$  can be derived by directly using the derivation method for each seller's profit function and then solving  $m$  simultaneous equations as above. However, for complicated function forms (e.g., more complicated loss function rather than the used one), since the number of sellers  $m$  can be quite large in practice, it may be difficult to derive analytical expressions by solving a large number of simultaneous equations each with complex forms. Specifically, the  $m$  equations are highly coupled, i.e., each with all  $\tau_i, i = 1, 2, \dots, m$ , and eliminating the similar terms to simplify the equations as we did in Eq. 12 is not always feasible. Therefore, we propose an approximate method which makes each equation with a single  $\tau_i$  and independent from others. Note that the approximate approach is proposed to deal with the case where direct derivation would fail rather than to improve the efficiency. Thus we take a different privacy loss function form for the sellers as an example where the direct derivation is not practically feasible in order to illustrate the mean-field method. Specifically, we replace Eq. 11 with  $\mathbf{L}_i(\tau_i) = \lambda_i \chi_i \tau_i^2$ .

The approximation is based on the mean-field theory [49], which deals with situations that involve a great number of agents, i.e, sellers in our context. When there are a great number of sellers in Nash game, it's reasonable to expect that a single seller has a *tiny* (infinitesimal) influence on the equilibrium and is affected by other sellers through a mean-field state, which we formulate as the weighted mean of all sellers' strategies,  $\bar{\tau}$ .

$$\bar{\tau} = \frac{\sum_{i=1}^m \omega_i \tau_i}{m}. \quad (16)$$

The mean-field state  $\bar{\tau}$  indicates the overall data fidelity provided by sellers at equilibrium and is not intensively affected by the data fidelity from one specific seller.



Using the new privacy loss function, the profit function of seller  $\mathcal{S}_i$  in Eq. 3 is changed into

$$\Psi_i(p^D, \tau_i) = p^D(\chi_i \tau_i) - \lambda_i \chi_i \tau_i^2. \quad (17)$$

Using  $\bar{\tau}$ ,  $\chi_i$  can be simplified as  $N \frac{\omega_i \tau_i}{m \bar{\tau}}$ . Since  $\bar{\tau}$  is not strongly affected by specific  $\tau_i$ , we can easily derive the first-order derivative of each seller's profit function  $\Psi_i(p^D, \tau_i)$  with respect to  $\tau_i$  and let them equal to zero.

$$\begin{cases} p^D \cdot N \frac{\omega_1 \tau_1^2}{m \bar{\tau}} - \lambda_1 \cdot N \frac{\omega_1 \tau_1^3}{m \bar{\tau}} = 0 \\ p^D \cdot N \frac{\omega_2 \tau_2^2}{m \bar{\tau}} - \lambda_2 \cdot N \frac{\omega_2 \tau_2^3}{m \bar{\tau}} = 0 \\ \vdots \\ p^D \cdot N \frac{\omega_m \tau_m^2}{m \bar{\tau}} - \lambda_m \cdot N \frac{\omega_m \tau_m^3}{m \bar{\tau}} = 0. \end{cases}$$

We derive  $\mathcal{S}_i$ 's optimal strategy

$$\tau_i^* = \frac{2p^D}{3\lambda_i}, i = 1, 2, \dots, m. \quad (18)$$

Note that the second-order derivative  $\frac{\partial^2 \Psi_i(p^D, \tau_i)}{\partial \tau_i^2} < 0$ , so these solutions can maximize each seller's profit.

**Error Analysis.** We use fixed  $\bar{\tau}$  to replace  $\sum_{i=1}^m \frac{\omega_i \tau_i}{m}$  when deriving the derivatives. Such replacement is an approximation and its error depends on the form of the profit function. In this part, we analyze the error bound of the approximated mean-field approach.

*Theorem 1:* The exact weighted mean of all sellers' strategies by the direct derivation is defined as  $\bar{\tau}^{DD}$ , and the approximated one by the mean-field method is  $\bar{\tau}^{MF}$ . The error is  $\bar{\tau}^{DD} - \bar{\tau}^{MF}$ . Consider the case that the privacy loss function is  $\mathbf{L}_i(\tau_i) = \lambda_i \chi_i \tau_i^2$ . When the number of sellers  $m$  is large and by scaling  $\omega_1, \omega_2, \dots, \omega_m$  such that  $\frac{\omega_i}{\lambda_i} \leq \frac{1}{p^D m^2}$ , we get

$$-\frac{1}{6m^2} < \bar{\tau}^{DD} - \bar{\tau}^{MF} < \frac{1}{m} - \frac{2}{3m^2}.$$

Note that what makes sense is the proportional relationship among  $\omega_i, i = 1, 2, \dots, m$ , allowing us to arbitrarily scale them.

*Proof 1:* We first calculate the upper bound of  $\bar{\tau}^{DD} - \bar{\tau}^{MF}$ . By applying direct derivation to Eq. 17, we can get

$$2p^D \sum_{j=1}^m \omega_j \tau_j - p^D \omega_i \tau_i = 3\lambda_i \tau_i \sum_{j=1}^m \omega_j \tau_j - \lambda_i \omega_i \tau_i^2.$$

By splitting  $\sum_{j=1}^m \omega_j \tau_j$  into  $\sum_{j=1, j \neq i}^m \omega_j \tau_j$  and  $\omega_i \tau_i$ , we obtain a quadratic equation about  $\tau_i$  by deforming the above formula, and using root formula for the quadratic equation, we can get

$$\tau_i^* = \frac{p^D \omega_i - 3\lambda_i \Sigma_{\tau-i} + \sqrt{(3\lambda_i \Sigma_{\tau-i} - p^D \omega_i)^2 + 16p^D \lambda_i \omega_i \Sigma_{\tau-i}}}{4\lambda_i \omega_i}, \quad (19)$$

where  $\Sigma_{\tau-i} = \sum_{j=1, j \neq i}^m \omega_j \tau_j$ . With the constraint  $\frac{\omega_i}{\lambda_i} \leq \frac{1}{p^D m^2}$ , we can justify that  $3\lambda_i \Sigma_{\tau-i} - p^D \omega_i > 0$  when  $m$  is very large. Thus according to  $\sqrt{x+y} < \sqrt{x} + \sqrt{y}$ , we can scale and deform the above formula to get

$$\omega_i \tau_i^* < \frac{\sqrt{16p^D \lambda_i \omega_i \Sigma_{\tau-i}}}{4\lambda_i}.$$

Further simplifying and scaling the above formula, we can get

$$\omega_i \tau_i^* < \sqrt{p^D \frac{\omega_i}{\lambda_i} \Sigma_{\tau-i}} \leq \sqrt{p^D \frac{\omega_i}{\lambda_i} \sum_{j=1}^m \omega_j \tau_j^*},$$

which applies to all  $\tau_i^*, i = 1, 2, \dots, m$ . Then, by adding  $m$  inequalities together and simplifying it, we can obtain

$$\sum_{i=1}^m \omega_i \tau_i^* < \left( \sum_{i=1}^m \sqrt{p^D \frac{\omega_i}{\lambda_i}} \right)^2,$$

and thus

$$\bar{\tau}^{DD} = \frac{1}{m} \sum_{i=1}^m \omega_i \tau_i^* < \frac{1}{m} \left( \sum_{i=1}^m \sqrt{p^D \frac{\omega_i}{\lambda_i}} \right)^2.$$

Additionally, using Eqs. 16 and 18, we can derive  $\bar{\tau}^{MF}$  as below.

$$\bar{\tau}^{MF} = \frac{1}{m} \sum_{i=1}^m \frac{2p^D \omega_i}{3\lambda_i}.$$

Then we use  $\frac{\omega_i}{\lambda_i} \leq \frac{1}{p^D m^2}$  and get

$$\begin{aligned} \bar{\tau}^{DD} - \bar{\tau}^{MF} &< \frac{1}{m} \left( \sum_{i=1}^m \sqrt{p^D \frac{\omega_i}{\lambda_i}} \right)^2 - \frac{1}{m} \sum_{i=1}^m \frac{2p^D \omega_i}{3\lambda_i} \\ &\leq \frac{1}{m} \left( \sum_{i=1}^m \sqrt{\frac{1}{m^2}} \right)^2 - \frac{1}{m} \sum_{i=1}^m \frac{2}{3m^2} \\ &= \frac{1}{m} - \frac{2}{3m^2}. \end{aligned}$$

Next, we calculate the lower bound. Since  $(3\lambda_i \Sigma_{\tau-i} - p^D \omega_i)^2 + 16p^D \lambda_i \omega_i \Sigma_{\tau-i} > (p^D \omega_i + 3\lambda_i \Sigma_{\tau-i})^2$ , using Eq. 19 we can get

$$\begin{aligned} \bar{\tau}^{DD} &= \frac{1}{m} \sum_{i=1}^m \omega_i \tau_i^* \\ &> \frac{1}{m} \sum_{i=1}^m \frac{p^D \omega_i - 3\lambda_i \Sigma_{\tau-i} + \sqrt{(p^D \omega_i + 3\lambda_i \Sigma_{\tau-i})^2}}{4\lambda_i} \\ &= \frac{1}{m} \sum_{i=1}^m \frac{p^D \omega_i}{2\lambda_i}, \end{aligned}$$

and using  $\frac{\omega_i}{\lambda_i} \leq \frac{1}{p^D m^2}$  again, we get

$$\begin{aligned} \bar{\tau}^{DD} - \bar{\tau}^{MF} &= \frac{1}{m} \sum_{i=1}^m \omega_i \tau_i^* - \frac{1}{m} \sum_{i=1}^m \frac{2p^D \omega_i}{3\lambda_i} \\ &> \frac{1}{m} \sum_{i=1}^m \frac{p^D \omega_i}{2\lambda_i} - \frac{1}{m} \sum_{i=1}^m \frac{2p^D \omega_i}{3\lambda_i} \geq -\frac{1}{6m^2}. \end{aligned}$$

Therefore, Theorem 1 holds.

Through the above error analysis, we draw the following empirical conclusion: by scaling the value of  $\omega_i$  ( $i = 1, 2, \dots, m$ ) to satisfy  $\frac{\omega_i}{\lambda_i} \leq \frac{1}{p^D m^2}$ , the error of the mean-field approximation method will be bounded in an acceptable range and decrease with increasing  $m$  when  $m$  is very large. When  $m$  approaches infinity, the error is approximately zero. This

result is in line with the mean-field theory [49]. When the number of sellers  $m$  is big, our proposed mean-field method appears reasonable in terms of error.

2) *Expression of  $p^{D*}$  in Stage 2:* We give the expression of  $p^{D*}$  for the broker by direct derivation.

**Direct Derivation.** By substituting Eq. 15 into Eq. 10, we get

$$\chi_i^* = N \frac{\omega_i \tau_i^*}{\sum_{j=1}^m \omega_j \tau_j^*} = N \frac{\sqrt{\frac{\omega_i}{\lambda_i}}}{\sum_{j=1}^m \sqrt{\frac{\omega_j}{\lambda_j}}}.$$

Then we get

$$q^{D*} = \sum_{i=1}^m \chi_i^* \tau_i^* = \sum_{i=1}^m \frac{p^D}{2\lambda_i}.$$

Since  $q^M = h(q^D, \nu)$  is positively correlated to  $q^D$  and  $\nu$ , we instantiate  $h(q^D, \nu)$  as  $q^D \nu$ . We get

$$q^{M*} = q^{D*} \nu = \frac{1}{2} \sum_{i=1}^m \frac{1}{\lambda_i} p^D \nu.$$

By substituting  $q^{D*}$  and  $q^{M*}$  into  $\mathcal{A}$ 's profit function in Eq. 2, we get

$$\Omega(p^M, p^D, \tau) = p^M \cdot \left( \frac{1}{2} \sum_{i=1}^m \frac{1}{\lambda_i} p^D \nu \right) - C(N, \nu) - p^D \cdot \left( \frac{1}{2} \sum_{i=1}^m \frac{1}{\lambda_i} p^D \right).$$

We derive the first-order derivative with respect to  $p^D$  and let it equal to 0.

$$\frac{\partial \Omega(p^M, p^D, \tau)}{\partial p^D} = \frac{1}{2} \sum_{i=1}^m \frac{1}{\lambda_i} \nu p^M - \sum_{i=1}^m \frac{1}{\lambda_i} p^D = 0.$$

We can thus get the expression of  $p^{D*}$

$$p^{D*} = \frac{\nu p^M}{2}. \quad (20)$$

Note that the second-order derivative  $\frac{\partial^2 \Omega(p^M, p^D, \tau)}{\partial p^{D^2}} < 0$ , so the solution can maximize the broker's profit.

3) *Value of  $p^{M*}$  in Stage 1:* We also use direct derivation in this stage, and by using the results in Sections IV-B1 and IV-B2, we can directly derive the value rather than the expression of  $p^{M*}$  for the buyer.

**Direct Derivation.** By substituting Eq. 15 and Eq. 20 into  $\mathcal{B}$ 's profit function in Eq. 1, we can obtain the profit of buyer  $\mathcal{B}$

$$\begin{aligned} \Phi(p^M, \tau) &= \theta_1 \ln(1 + \rho_1 q^{D*}) + \theta_2 \ln(1 + \rho_2 \nu) - p^M q^{M*} \\ &= \theta_1 \ln(1 + c_1 p^M) + \theta_2 \ln(1 + \rho_2 \nu) - \frac{c_2 \theta_1}{2} p^{M^2}, \end{aligned}$$

where  $c_1 = \frac{\rho_1 \nu}{4} \sum_{i=1}^m \frac{1}{\lambda_i}$  and  $c_2 = \frac{\nu^2}{2\theta_1} \sum_{i=1}^m \frac{1}{\lambda_i}$ . Then, we derive the first-order derivative of  $\Phi(p^M, \tau)$  as follows.

$$\frac{\partial \Phi(p^M, \tau)}{\partial p^M} = \frac{\theta_1 c_1}{1 + c_1 p^M} - c_2 \theta_1 p^M. \quad (21)$$

By letting  $\frac{\partial \Phi(p^M, \tau)}{\partial p^M}$  in Eq. 21 equal to zero, we obtain

$$c_1 c_2 \cdot p^{M^2} + c_2 \cdot p^M - c_1 = 0.$$

Using the characteristic root method, we find buyer  $\mathcal{B}$ 's optimal strategy  $p^{M*}$  (after discarding the negative solution).

$$p^{M*} = \frac{-c_2 + \sqrt{c_2^2 + 4c_1^2 c_2}}{2c_1 c_2}. \quad (22)$$

We justify that the second-order derivative  $\frac{\partial^2 \Phi(p^M, \tau)}{\partial p^{M^2}} = -\frac{\theta_1 c_1^2}{(1 + c_1 p^M)^2} - \theta_1 c_2 < 0$ , so the solution can maximize the buyer's profit.

Getting  $p^{M*}$ , we can determine the optimal value of  $p^{D*}$  by substituting  $p^{M*}$  into Eq. 20 and each seller's optimal value of  $\tau_i^*$  by substituting  $p^{D*}$  into Eq. 15. Till now, the complete optimal strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$  has been determined, based on which the market transaction can be conducted.

4) *Equilibrium Analysis:* We prove the existence and uniqueness of SNE in *Share*.

**Theorem 2:** The complete optimal strategy profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$  determined by backward induction approach uniquely constitutes SNE.

*Proof 2:* We prove the existence and uniqueness of SNE.

**Existence.** For the buyer, when the broker and sellers hold the optimal strategies in Eq. 20 and Eq. 15, the buyer's profit  $\Phi(p^M, \tau^*)$  only changes with  $p^M$ . In the process of deriving optimal  $p^{M*}$ , the first-order derivation is set to be 0 and the second-order is strictly less than 0, which means that the maximum profit is obtained at  $p^{M*}$ . Thus, Eq. 4 holds at  $p^{M*}$ . For the broker, when the buyer and sellers hold the optimal strategies,  $p^{M*}$  in Eq. 22 and  $\tau^*$  in Eq. 15, the broker's optimal profit can be obtained at  $p^{D*}$  since the profit function is strictly concave and has a single extreme point  $p^{D*}$ . Thus, Eq. 5 holds at  $p^{D*}$ . For the sellers, each seller determines each  $\tau_i^*$  simultaneously in the same way, and  $\tau_i^*, i = 1, 2, \dots, m$  are jointly decided. For each seller  $\mathcal{S}_i$ , her optimal strategy  $\tau_i^*$  is determined by letting the first-order derivation equal to 0. Since the profit function is strictly concave, the extreme point  $\tau_i^*$  maximizes seller  $\mathcal{S}_i$ 's profit if  $\tau_i^* \leq 1$ . Otherwise, when the extreme point is larger than 1, the optimal value  $\tau_i^* = 1$  can also maximize  $\mathcal{S}_i$ 's profit since the profit function is monotonically increasing in the feasible range of  $\tau_i$  and maximized at the right endpoint 1. Thus, Eq. 6 holds at  $\tau_i^*$ . Therefore, it's proved that SNE exists in our mechanism.

**Uniqueness.** For the buyer, since her profit function is strictly concave, the maximum profit is obtained only at  $p^{M*}$ . Any other value of  $p^M \neq p^{M*}$  will yield an inferior profit. Such result can be also explained by Convex Optimization [50], i.e., the strategy space of  $p^M$  is a convex and compact subspace of Euclidean space, and the profit function  $\Phi(\cdot)$  is a convex function of  $p^M$ , leading to the unique optimal  $p^{M*}$  that maximizes  $\Phi(\cdot)$ . Thus, Eq. 4 holds only at  $p^{M*}$ . For the broker, her profit function is also strictly concave and only has a single extreme point  $p^{D*}$ . Any other value of  $p^D \neq p^{D*}$  will lower the broker's profit. Thus, Eq. 5 holds only at  $p^{D*}$ .

For sellers, seller  $S_i$  can only have lower profit by deciding  $\forall \tau_i \neq \tau_i^*$ . If  $\tau_i^* \leq 1$ , seller  $S_i$ 's profit function is concave and only maximized at  $\tau_i^*$ . Otherwise, the profit can only be maximized at the right endpoint,  $\tau_i^* = 1$ , since the profit function is monotonically increasing. For each seller, if she chooses other  $\forall \tau_i \neq \tau_i^*$ , she can only have lower profit with  $\tau_{-i}^*$  kept the same. Thus, the unique Nash equilibrium among sellers is achieved and Eq. 6 holds only at  $\tau_i^*$ . Therefore, it's proved that other strategy profiles except our solution cannot satisfy SNE, which indicates the uniqueness of SNE in our mechanism.

### C. Complete Data Trading Dynamics

We summarize the complete dynamics of data markets in Algorithm 1, which integrates the equilibrium solving process in Section IV-B.

The first phase is *Initialization*. We assume that each party can report their specific input parameters which can be fitted based on the historical experiences. The buyer sets appropriate parameters  $\theta_1, \theta_2, \rho_1, \rho_2$  for her utility function and proposes queries  $\mathbf{Q}$  with performance parameter  $\nu$  which need to be solved by the product (Line 2). Note that the product is not restricted in forms decided by the broker while the access channel for data buyers is uniformly set as queries. The broker crystallizes the size  $N$  of data she can handle, determines  $\sigma_k, k \in \{0, 1, 2, 3, 4, 5\}$  for her cost function, and maintains the weights  $\omega_i, i = 1, 2, \dots, m$  of sellers' datasets (Line 3). To decide the real weights before the first transaction, the broker can use dummy buyers to iterate several times where Shapley value can be used to evaluate the sellers' datasets. Sellers give their privacy sensitivity  $\lambda_i, i = 1, 2, \dots, m$  (Line 4).

The second phase is *Strategy Decision*. Using the strategy mechanism, buyer  $\mathcal{B}$ , broker  $\mathcal{A}$ , and each seller  $S_i$  give product price  $p^{M*}$ , data price  $p^{D*}$ , and data fidelity  $\tau_i^*$  in order according to Eqs. 22, 20, 15, respectively (Line 6).

Then *Data Transaction* between the broker and sellers begins. The data quantity chosen from each seller can be calculated according to Eq. 10 (Line 8). Each seller randomly picks  $\chi_i^*$ -sized dataset (Line 10) and pre-processes it for privacy protection based on  $\epsilon_i^*$  calculated from Eq. 9 (Lines 11-12). After that, seller  $S_i$  gives her protected dataset  $D_i^t$  to the broker in exchange for compensation  $\mathbf{PS}_i^*$  (Line 13).

The next phase is *Product Production*. The broker collects the data as  $D^t$  and uses it to make the product (Line 15). Moreover, the weights of sellers' datasets are updated by the broker based on their corresponding contributions to the data product (Line 16). We give one update formula based on Shapley value as an example:  $\omega_i' = 0.2\omega_i + 0.8S\mathcal{V}_i$ , where  $S\mathcal{V}_i$  is the Shapley value of  $D_i^t$  to the product and the coefficient 0.8 indicates to what extent the historical performance of data can be useful for the current task. The updated weights  $\omega_i'$  can be used in the subsequent transaction.

The last phase is *Product Transaction* between the broker and the buyer. The broker gives the product (strictly speaking, the query answers based on the product) to the buyer and the buyer pays  $\mathbf{PB}^*$  to the broker (Line 18). So far, the current

---

### Algorithm 1: Data trading dynamics.

---

```

1 %% Initialization;
2 From the current buyer  $\mathcal{B}$ , demanded queries  $\mathbf{Q}$  and parameters
    $\nu, \theta_1, \theta_2, \rho_1, \rho_2$  are provided;
3 From broker  $\mathcal{A}$ ,  $N, \sigma_k (k \in \{0, 1, 2, 3, 4, 5\}), \omega_i (i = 1, 2, \dots, m)$ 
   are given;
4 From existing  $m$  sellers, each seller  $S_i$  decides  $\lambda_i$ ;
5 %% Strategy Decision;
6 Through three-stage Stackelberg-Nash game, the optimal strategy
   profile  $\langle p^{M*}, p^{D*}, \tau^* \rangle$  is determined by the buyer, the broker,
   and sellers, respectively;
7 %% Data Transaction;
8 The quantity of data each seller can sell,  $\chi^*$ , is calculated according
   to Eq. 10;
9 for each seller  $S_i, i = 1, 2, \dots, m$  do
10   Randomly pick  $\chi_i^*$  data pieces from her dataset  $D_i$ ;
11   Calculate  $\epsilon_i^*$  from the strategy  $\tau_i^*$  according to Eq. 9;
12   Conduct LDP with  $\epsilon_i^*$  on her  $\chi_i^*$ -sized dataset, and then give
       the protected  $D_i^t$  to broker  $\mathcal{A}$ ;
13 Broker  $\mathcal{A}$  gets data from sellers to form dataset  $D^t$  for production
   and pays compensation  $\mathbf{PS}_i^*$  to each seller;
14 %% Product Production;
15 Broker  $\mathcal{A}$  then uses  $D^t$  to produce the data product as well as
   computes the answers to queries  $\mathbf{Q}$ ;
16 After manufacturing the product, broker  $\mathcal{A}$  updates  $\omega_1, \omega_2, \dots, \omega_m$ 
   (might scale down or normalized as needed) based on the
   contribution to the product from each seller's  $D_i^t$ ;
17 %% Product Transaction;
18 Broker  $\mathcal{A}$  gives the product to buyer  $\mathcal{B}$  (by returning the query
   answers), and meantime buyer  $\mathcal{B}$  pays  $\mathbf{PB}^*$  to broker  $\mathcal{A}$ .

```

---

data transaction among buyer  $\mathcal{B}$ , broker  $\mathcal{A}$ , and sellers  $S_i, i = 1, 2, \dots, m$  has finished. When the next buyer comes, the next transaction will start.

## V. EXPERIMENTS

In this section, we present experimental studies validating the effectiveness and efficiency of *Share*. We first describe our experiment setup including the datasets, parameter settings, and measurement metrics in Section V-A. Sections V-B and V-C show the results verifying the effectiveness and efficiency of *Share*, respectively. Section V-D shows the effects of the main parameters used in *Share*.

### A. Experiment Setup

We conduct experiments on a machine with an Intel Core i7-11700KF running Ubuntu with 64GB memory. The Shapley value is calculated based on Monte Carlo Method [51]. Laplace mechanism [52], a technique for achieving LDP, is applied to each record to adjust data fidelity for each seller.

**Datasets.** We use a real dataset, Combined Cycle Power Plant (CCPP) [53], which contains 9,568 data points with four features. The buyer's demanded query task is to get the prediction of net hourly electrical energy output given the specific conditions of features, which can be understood as the analytical query extensively used in analytical database. A linear regression model is considered as the data product which the broker manufactures to serve the queries, and explained variance is used to measure the performance. We randomly choose a training dataset (the data of sellers) with a size of 9,000, and the 568 data records left are used for validation

(based on which the queries of the buyer are generated). In the real world, the datasets of sellers can be the same in quality (which makes it easy to randomly choose sellers to buy data), or vary in quality, which is the case we deal with in *Share*. To stimulate the distinction in data quality, we first sort data by quality measured by Shapley value, which indicates the contribution of each data record to regression. Then by distributing data in decreasing quality over sellers, each seller owns 90 data records with different quality. Besides the real dataset, we augment CCPP through replication and Gaussian noise  $\mathcal{N}(0, 0.1^2)$  injection to generate a synthetic dataset with a size of 1,000,000 to test the efficiency of *Share*.

**Parameter Settings.** Our parameters include the number of sellers  $m$ , the total data quantity  $N$ , the required explained variance  $\nu$ , the individual parameters of each party's profit (i.e., buyer  $\mathcal{B}$ 's  $\theta_1, \theta_2, \rho_1, \rho_2$  related to utility, broker  $\mathcal{A}$ 's cost parameters  $\sigma_k, k \in \{0, 1, 2, 3, 4, 5\}$ , and seller  $\mathcal{S}_i$ 's privacy sensitivity  $\lambda_i, i = 1, 2, \dots, m$ ), and the initial weights  $\omega_i, i = 1, 2, \dots, m$  of sellers. We set  $m = 100$ ,  $N = 500$ , and  $\nu = 0.8$ . The utility parameters of the buyer are set as  $\theta_1 = 0.5, \theta_2 = 0.5, \rho_1 = 0.5, \rho_2 = 250$  (in order to balance the impacts of product performance and dataset quality). The cost parameters of the broker are related to the practical manufacturing situation and are set as default values  $\sigma_0 = 1 \times 10^{-3}, \sigma_1 = -2, \sigma_2 = -3, \sigma_3 = 1 \times 10^{-3}, \sigma_4 = 2 \times 10^{-3}, \sigma_5 = 1 \times 10^{-3}$ . Sellers'  $\lambda_i, i = 1, 2, \dots, m$  are picked randomly in  $(0, 1)$ .  $\omega_1, \omega_2, \dots, \omega_m$  are initially generated by using a dummy buyer to iterate the mechanism which takes five times to stabilize the profits. We consider buyer  $\mathcal{B}$  as a general buyer coming after several transactions have finished. Shapley values of sellers' datasets can be calculated after regression to update the weights for the next transaction.

**Measurement Metrics.** Four main indexes are adopted to evaluate the mechanism concerning both effectiveness and efficiency. We will show the results of using direct derivation for equilibrium solving. The mean-field approach (used when direct derivation fails) performs the same on the metrics.

- *Profit.* The profits of the buyer, the broker, and sellers.
- *Social welfare.* The effect on the overall benefit, notated as a function of  $\tau$ , i.e.,  $\mathbf{SW}(\tau) = \mathbf{U}(\tau) - \sum_{i=1}^m \mathbf{L}_i(\tau_i) - \mathbf{C}$ . Note that the social welfare measures the total profits, correlated to the collected data (more precisely, data fidelity  $\tau$ ) yet independent of payments (prices) which only circulate among participants. The optimum social welfare can be derived by solving the social welfare maximization problem  $\max_{\tau} \mathbf{SW}$  and used to measure the social welfare level of our mechanism which is represented as the ratio to the social optimum.
- *Product quality.* The quality of the data product manufactured by the collected data with distinctive fidelity.
- *Runtime.* The time cost of executing the mechanism which evaluates the efficiency.

## B. Effectiveness

We first unilaterally change the strategies of the buyer, the broker, and sellers respectively to verify the profit maximization and the implied equilibrium. The social welfare and product quality are then evaluated with different numbers of sellers. We also investigate the effect of the inner Nash game.

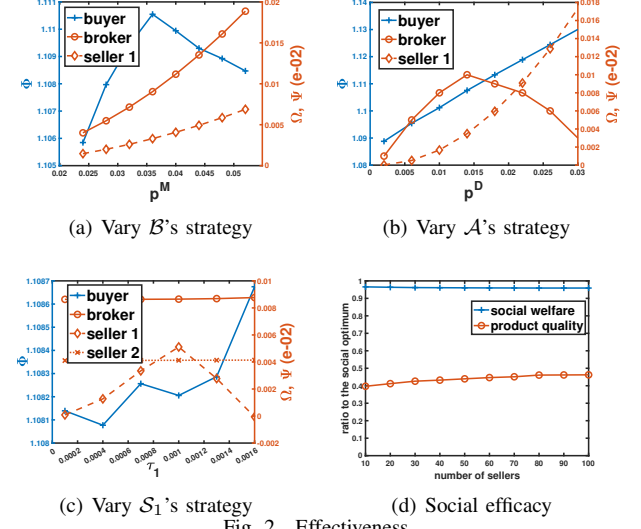


Fig. 2. Effectiveness.

Fig.2(a) shows the results of profits when we change  $p^M$  around the optimal strategy  $p^{M*}$  while maintaining the rest. Seller  $\mathcal{S}_1$  acts as a representative of sellers. It's found that the peak of the buyer's profit  $\Phi(\cdot)$  appears when her optimal strategy  $p^{M*} = 0.036$  determined in SNE is adopted (the monetary unit can adjust with how the utility/cost function is mapped into money). Whatever strategy the buyer chooses except  $p^{M*}$ , she will get a lower profit when all other participants' strategies are fixed. The change of the profits of the broker and the seller is intuitive. Specifically, with growing  $p^M$ , the broker can gain more profit, which can further add the compensations for sellers and make their profits higher.

Fig.2(b) shows the results of profits when we change  $p^D$  around the optimal strategy  $p^{D*} = 0.014$  while maintaining the rest. Similarly, it is found that the broker cannot increase her profit by unilaterally changing her strategy. The change of the profits of the buyer and seller is also intuitive. Specifically, the growing  $p^D$  brings more compensations to sellers, adding their profits. Due to more compensations, the dataset quality from sellers can therefore be improved, which causes the rise of the buyer's profit.

Fig.2(c) shows the results of profits when we change  $\tau_1$  around the optimal strategy  $\tau_1^* = 0.001$  while maintaining the rest. The first two sellers  $\mathcal{S}_1$  and  $\mathcal{S}_2$  are chosen as representatives. It is the same that the seller who changes her strategy unilaterally gets no more profit. Even if one seller changes her strategy, the broker can nearly keep her profit as before, benefiting from the *inner* competition among sellers which is formulated as a Nash game. Specifically, the effect of sellers' bounded rationality is almost limited among sellers and is corrected automatically in Stage 3, which signifies the

TABLE II  
COMPARING NASH WITH OTHERS.

	Nash	Random	Average
$q^M$	6.379658	2.013382	2.022149
$\Phi$	3.255229	2.995894	2.996831
$\Omega$	0.099928	0.031531	0.031669
$\Psi$	0.000500	0.000154	0.000161
SW	96%	90%	90%

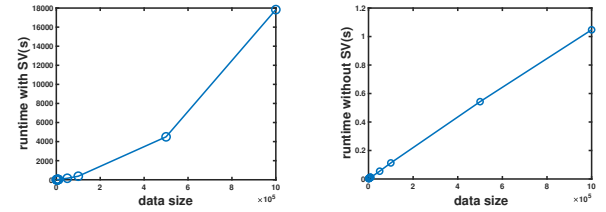
transparency of Stage 3 to the upper stages. The change of the buyer's profit may be due to the effect of the data on the regression, which is not always predictable, causing the irregular curve of  $\Phi(\cdot)$ . In theory, varying  $\tau_1$  surely makes differences on other sellers. However, since the number of sellers is large, this effect is *diluted* and negligible, making the profit of  $S_2$  almost unchanged.

Fig.2(d) shows the ratio results of social welfare and data product quality, collectively referred as social efficacy, compared to the optimum ones derived from the social welfare maximization problem. The proposed mechanism can achieve extremely high (over 95%) social welfare. As the number of sellers rises up, the social welfare slightly decreases, which implies that more strategic gaming among participants exacerbates the social inefficiency in terms of the overall profits. The product quality performs inferior to the social optimal result due to the selfish profit-seeking behaviors of participants, which, however, would still outperform the baselines as shown in the following justification of Nash game. Better product can be acquired when more sellers (thus, more data) engage, implying the significance of data circulation.

**Inner Nash Game.** To verify the effectiveness of using Nash game to formulate Stage 3, we implement the mechanism compared to the baselines of using Random and Average strategies to select data. Table II shows the product quality  $q^M$ , the profit results of the buyer  $\Phi$ , the broker  $\Omega$ , and the average level of sellers  $\Psi$ , as well as the social welfare SW (represented as the ratio to the optimum) through Nash game, Random, and Average respectively with the parameters kept the same. We can observe that the Nash-based seller selection outperforms the baselines in terms of the data product quality, the profits of all three parties, and the social welfare. The product quality based on the data selected by Nash game is the highest, which explicitly shows the effectiveness of the seller selection results. In terms of individual profits, not only data sellers benefit from their inner benign competition, the profits of both the buyer and the broker also increase, which indicates the advantage of the inner Nash game to the upper stages. In terms of social benefits, the Nash game modelling achieves the highest social welfare compared to the optimal one, implying its effectiveness in data markets both individually and collectively.

Note that the buyer's profit is much more than the broker's and sellers', which is in line with the property of Stackelberg game (in favor of the leader) and consistent with the desired effect in demand-driven markets. The buyer as the transaction initiator can create value using the demanded product and gain long-term benefits (e.g., the huge revenue Ford earns owing to

the business decision based on the acquired query answers), while the broker or the sellers make profit from the one-shot transaction which is relatively lower.



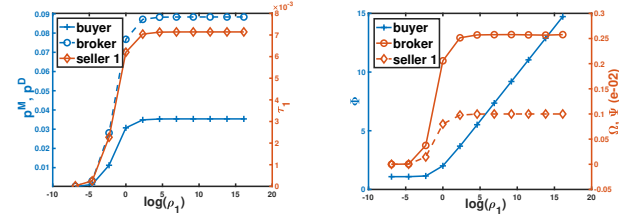
(a) Runtime with SV (b) Runtime without SV  
Fig. 3. Efficiency.

### C. Efficiency

Fig.3(a) and Fig.3(b) show the runtime of the proposed data trading algorithm with and without Shapley value to update weights. We use the synthetic dataset with 1,000,000 data records and adjust the number of sellers  $m$  from 5 to 10,000 while fixing the other parameters and the average number of data records chosen from each seller as 100. Fig.3(a) shows that the runtime grows as  $m$  goes higher but with an acceptable rate. Even when  $m = 10,000$ , it does not take too much time. While our mechanism contains a time-consuming part to calculate Shapley values, Fig.3(b) shows that our mechanism without Shapley value calculation can run very fast with a linear time complexity, which corresponds to the theoretical analysis of the complexity of Algorithm 1.

### D. Parameter Influence

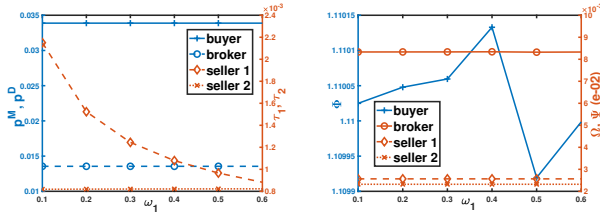
In this section, we make sensitivity analyses on the major parameters in our mechanism and investigate how the parameters affect the strategies and profits of the three parties.



(a) Strategy vs.  $\rho_1$  (b) Profit vs.  $\rho_1$   
Fig. 4. Effect of  $\rho_1$ .

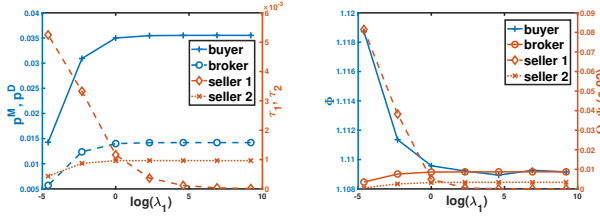
Fig.4(a) and Fig.4(b) present the effect of  $\rho_1$  on strategies and profits. Note that  $\rho_1$  is a parameter relevant to the buyer's sensitivity to dataset quality, which objectively reflects the relationship between dataset quality and product utility. Fig.4(a) shows that too small  $\rho_1$  can hardly lead to effective markets because of the buyer's indifference on the data. When  $\rho_1$  reaches a certain level, all the strategies stay the same and the market reaches equilibrium. The influence of  $\rho_1$  is limited within the utility for the buyer and can no longer disturb the market equilibrium, which may be due to common sense that the dataset quality can't increase unlimitedly and even with sharper sensitivity to the data, higher prices wouldn't bring about better data anymore. Fig.4(b) shows that the profit of the buyer surges as  $\rho_1$  increases because she will get more utility from the raise of dataset quality. When  $\rho_1$  is big enough, the increase of  $\rho_1$  has little effect on the profits of the broker and sellers, which can be explained by the trend of strategies.





(a) Strategy vs.  $\omega_1$  (b) Profit vs.  $\omega_1$   
Fig. 5. Effect of  $\omega_1$ .

Fig.5(a) and Fig.5(b) present the effect of  $\omega_1$  on strategies and profits. Note that  $\omega_1, \omega_2, \dots, \omega_m$  are the weights of sellers' datasets and assess the sellers' data in previous transactions. We select  $\mathcal{S}_1$  and  $\mathcal{S}_2$  as representatives. Fig.5(a) shows that  $\omega_1$  only affects the strategy of the corresponding seller  $\mathcal{S}_1$ . The strategies of the buyer and the broker remain the same because  $\omega_1$  only affects the inner-seller competition. Since the number of sellers is large, varying  $\omega_1$  makes little difference on other sellers, making the strategy of seller  $\mathcal{S}_2$  almost unchanged. Fig.5(b) shows that when  $\omega_1$  varies from 0.1 to 0.6, all profits except the buyer's are stable. Once  $\omega_1$  gets an inappropriate value, the data of seller  $\mathcal{S}_1$  won't work as expected and affects the profit of the buyer, leading to the unsmooth curve of  $\Phi(\cdot)$ .



(a) Strategy vs.  $\lambda_1$  (b) Profit vs.  $\lambda_1$   
Fig. 6. Effect of  $\lambda_1$ .

Fig.6(a) and Fig.6(b) show the effect of  $\mathcal{S}_1$ 's parameter  $\lambda_1$  on strategies and profits. Note that  $\lambda_i$  is related to seller  $\mathcal{S}_i$ 's privacy sensitivity. Fig.6(a) shows that  $\tau_1$  sinks with increasing  $\lambda_1$  since seller  $\mathcal{S}_1$  will strengthen her data protection if more sensitive to privacy risks.  $p^M$  and  $p^D$  increase possibly because higher prices may be provided to encourage conservative sellers to offer high-fidelity data in spite of heavy privacy risks. Fig.6(b) shows that  $\lambda_1$  mainly influences the buyer's and the corresponding seller  $\mathcal{S}_1$ 's profits. The profit of  $\mathcal{S}_1$  decreases because bigger  $\lambda_1$ , more privacy loss  $\mathcal{S}_1$  will suffer. The profit of buyer  $\mathcal{B}$  dives probably because the seller would enhance the protection on her data when faced with huge privacy risks, thus lowering the data fidelity and further harming the buyer's profit. The profit of the broker remains unchanged because the broker herself does not rely on the data fidelity but just transfers the data from the sellers to the buyer.

## VI. CONCLUSION AND FUTURE WORK

We present *Share*, the first incentive data market framework applicable to demand-driven scenarios based on a three-stage Stackelberg-Nash game. The profit maximization for all participants and the *buyer-broker-sellers* market flow are fulfilled by considering the mutual interaction among three parties as a three-stage Stackelberg game, in which the absolute pricing for

data is also realized. We address the seller selection problem by considering the inter-seller competition as a Nash game. To derive the Stackelberg-Nash Equilibrium, backward induction is used, and a novel mean-field approximation with provable guarantees is proposed. Our proposed data market framework performs well on real and synthetic datasets in terms of both effectiveness and efficiency.

There are some interesting issues in the framework remained to be perfected, e.g., how to solve the challenge of parameter fitting for each party with deficient historical experiences, and how to encourage participants to honestly report those parameters via mechanisms rather than market regulations. Besides, there are other practical scenarios in data markets worth to be researched into, e.g., how to formulate a market with multiple buyers or streaming buyers, and how to support supply-driven data transactions. Nevertheless, we introduce a game theory based framework catering to a typical scenario with some of the many questions investigated, which has practical significance and may take step forward in data market research direction.

## REFERENCES

- [1] "Boston database meeting," 2023, [https://www.linkedin.com/posts/seemohan\\_45-worldwide-database-researchers-brainstorm-activity-7121547469573824512-hsB-?](https://www.linkedin.com/posts/seemohan_45-worldwide-database-researchers-brainstorm-activity-7121547469573824512-hsB-?)
- [2] J. Pei, R. C. Fernandez, and X. Yu, "Data and AI model markets: Opportunities for data and model sharing, discovery, and integration," *Proc. VLDB Endow.*, vol. 16, no. 12, pp. 3872–3873, 2023. [Online]. Available: <https://www.vldb.org/pvldb/vol16/p3872-pei.pdf>
- [3] R. C. Fernandez, P. Subramaniam, and M. J. Franklin, "Data market platforms: Trading data assets to solve data problems," *Proc. VLDB Endow.*, vol. 13, no. 11, pp. 1933–1947, 2020. [Online]. Available: <http://www.vldb.org/pvldb/vol13/p1933-fernandez.pdf>
- [4] A. Agarwal, M. Dahleh, and T. Sarkar, "A marketplace for data: An algorithmic solution," in *Proceedings of the 2019 ACM Conference on Economics and Computation*. ACM, 2019, pp. 701–726.
- [5] L. Chen, P. Koutris, and A. Kumar, "Towards model-based pricing for machine learning in a data marketplace," in *SIGMOD*, P. A. Boncz, S. Manegold, A. Ailamaki, A. Deshpande, and T. Kraska, Eds. ACM, 2019, pp. 1535–1552. [Online]. Available: <https://doi.org/10.1145/3299869.3300078>
- [6] J. Liu, J. Lou, J. Liu, L. Xiong, J. Pei, and J. Sun, "Dealer: An end-to-end model marketplace with differential privacy," *Proc. VLDB Endow.*, vol. 14, no. 6, pp. 957–969, 2021. [Online]. Available: <http://www.vldb.org/pvldb/vol14/p957-liu.pdf>
- [7] A. Tanner, "How data brokers make money off your medical records," *Scientific American*, vol. 314, no. 2, pp. 26–29, 2016.
- [8] P. Koutris, P. Upadhyaya, M. Balazinska, B. Howe, and D. Suciu, "Query-based data pricing," in *PODS*. ACM, 2012, pp. 167–178.
- [9] R. Jia, D. Dao, B. Wang, F. A. Hubis, N. M. Gurel, B. Li, C. Zhang, C. Spanos, and D. Song, "Efficient task-specific data valuation for nearest neighbor algorithms," *Proceedings of the VLDB Endowment*, vol. 12, no. 11, pp. 1610–1623, 2019.
- [10] Y. Li, X. Yu, and N. Koudas, "Data acquisition for improving machine learning models," *Proc. VLDB Endow.*, vol. 14, no. 10, pp. 1832–1844, 2021. [Online]. Available: <http://www.vldb.org/pvldb/vol14/p1832-li.pdf>
- [11] X. Cao, Y. Chen, and K. J. R. Liu, "Data trading with multiple owners, collectors, and users: An iterative auction mechanism," *IEEE Trans. Signal Inf. Process. over Networks*, vol. 3, no. 2, pp. 268–281, 2017. [Online]. Available: <https://doi.org/10.1109/TSIPN.2017.2668144>
- [12] R. C. Fernandez, "Protecting data markets from strategic buyers," in *SIGMOD '22: International Conference on Management of Data*, Philadelphia, PA, USA, June 12 - 17, 2022, Z. Ives, A. Bonifati, and A. E. Abbadi, Eds. ACM, 2022, pp. 1755–1769. [Online]. Available: <https://doi.org/10.1145/3514221.3517855>
- [13] L. S. Shapley, "A value for n-person games," *Contributions to the Theory of Games*, vol. 2, no. 28, pp. 307–317, 1953.

- [14] H. Von Stackelberg, *Market structure and equilibrium*. Springer Science & Business Media, 2010.
- [15] B. An, M. Xiao, A. Liu, X. Xie, and X. Zhou, "Crowdsensing data trading based on combinatorial multi-armed bandit and stackelberg game," in *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. IEEE, 2021, pp. 253–264.
- [16] J. C. Duchi, M. I. Jordan, and M. J. Wainwright, "Local privacy and statistical minimax rates," in *54th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2013, 26-29 October, 2013, Berkeley, CA, USA*. IEEE Computer Society, 2013, pp. 429–438. [Online]. Available: <https://doi.org/10.1109/FOCS.2013.53>
- [17] J. F. Nash Jr, "Equilibrium points in n-person games," *Proceedings of the national academy of sciences*, vol. 36, no. 1, pp. 48–49, 1950.
- [18] Bloomberg, "https://www.bloomberg.com/professional/product/market-data/," 1981.
- [19] DAWEX, "https://www.dawex.com/en/," 2015.
- [20] P. Koutiris, P. Upadhyaya, M. Balazinska, B. Howe, and D. Suciu, "Toward practical query pricing with querymarket," in *SIGMOD*. ACM, 2013, pp. 613–624.
- [21] M. Lei, X. Zhang, L. Chu, Z. Wang, P. S. Yu, and B. Fang, "Finding route hotspots in large labeled networks," *IEEE Trans. Knowl. Data Eng.*, vol. 33, no. 6, pp. 2479–2492, 2021. [Online]. Available: <https://doi.org/10.1109/TKDE.2019.2956924>
- [22] Y. Xu, C. Ma, Y. Fang, and Z. Bao, "Efficient and effective algorithms for generalized densest subgraph discovery," *Proc. ACM Manag. Data*, vol. 1, no. 2, pp. 169:1–169:27, 2023. [Online]. Available: <https://doi.org/10.1145/3589314>
- [23] K. Huang, H. Hu, Q. Ye, K. Tian, B. Zheng, and X. Zhou, "TED: towards discovering top-k edge-diversified patterns in a graph database," *Proc. ACM Manag. Data*, vol. 1, no. 1, pp. 51:1–51:26, 2023. [Online]. Available: <https://doi.org/10.1145/3588736>
- [24] M. Chen, Y. Zhao, Y. Liu, X. Yu, and K. Zheng, "Modeling spatial trajectories with attribute representation learning (extended abstract)," in *39th IEEE International Conference on Data Engineering, ICDE 2023, Anaheim, CA, USA, April 3-7, 2023*. IEEE, 2023, pp. 3813–3814. [Online]. Available: <https://doi.org/10.1109/ICDE55515.2023.00333>
- [25] H. Xie, Y. Fang, Y. Xia, W. Luo, and C. Ma, "On querying connected components in large temporal graphs," *Proc. ACM Manag. Data*, vol. 1, no. 2, pp. 170:1–170:27, 2023. [Online]. Available: <https://doi.org/10.1145/3589315>
- [26] L. Chu, Z. Wang, J. Pei, Y. Zhang, Y. Yang, and E. Chen, "Finding theme communities from database networks," *Proc. VLDB Endow.*, vol. 12, no. 10, pp. 1071–1084, 2019. [Online]. Available: <http://www.vldb.org/pvldb/vol12/p1071-chu.pdf>
- [27] J. Wang and Q. Zhang, "Disaggregated database systems," in *Companion of the 2023 International Conference on Management of Data, SIGMOD/PODS 2023, Seattle, WA, USA, June 18-23, 2023*, S. Das, I. Pandis, K. S. Candan, and S. Amer-Yahia, Eds. ACM, 2023, pp. 37–44. [Online]. Available: <https://doi.org/10.1145/3555041.3589403>
- [28] R. Wang, J. Wang, P. Kadam, M. T. Özsu, and W. G. Aref, "dlsm: An lsm-based index for memory disaggregation," in *39th IEEE International Conference on Data Engineering, ICDE 2023, Anaheim, CA, USA, April 3-7, 2023*. IEEE, 2023, pp. 2835–2849. [Online]. Available: <https://doi.org/10.1109/ICDE55515.2023.00217>
- [29] Y. Zhang, Q. Ye, R. Chen, H. Hu, and Q. Han, "Trajectory data collection with local differential privacy," *Proc. VLDB Endow.*, vol. 16, no. 10, pp. 2591–2604, 2023. [Online]. Available: <https://www.vldb.org/pvldb/vol16/p2591-chen.pdf>
- [30] R. Du, Q. Ye, Y. Fu, H. Hu, J. Li, C. Fang, and J. Shi, "Differential aggregation against general colluding attackers," in *39th IEEE International Conference on Data Engineering, ICDE 2023, Anaheim, CA, USA, April 3-7, 2023*. IEEE, 2023, pp. 2180–2193. [Online]. Available: <https://doi.org/10.1109/ICDE55515.2023.00169>
- [31] Z. Cong, X. Luo, J. Pei, F. Zhu, and Y. Zhang, "Data pricing in machine learning pipelines," *Knowl. Inf. Syst.*, vol. 64, no. 6, pp. 1417–1455, 2022. [Online]. Available: <https://doi.org/10.1007/s10115-022-01679-4>
- [32] J. Pei, "A survey on data pricing: from economics to data science," *IEEE Trans. Knowl. Data Eng.*, 2021.
- [33] J. Pei, F. Zhu, Z. Cong, X. Luo, H. Liu, and X. Mu, "Data pricing and data asset governance in the AI era," in *KDD '21: The 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, Singapore, August 14-18, 2021*, F. Zhu, B. C. Ooi, and C. Miao, Eds. ACM, 2021, pp. 4058–4059. [Online]. Available: <https://doi.org/10.1145/3447548.3470818>
- [34] R. P. McAfee and J. McMillan, "Auctions and bidding," *Journal of economic literature*, vol. 25, no. 2, pp. 699–738, 1987.
- [35] A. Nagurney and P. Dutta, "Supply chain network competition among blood service organizations: a generalized nash equilibrium framework," *Ann. Oper. Res.*, vol. 275, no. 2, pp. 551–586, 2019. [Online]. Available: <https://doi.org/10.1007/s10479-018-3029-2>
- [36] Y. Zhao, K. Zheng, J. Guo, B. Yang, T. B. Pedersen, and C. S. Jensen, "Fairness-aware task assignment in spatial crowdsourcing: Game-theoretic approaches," in *37th IEEE International Conference on Data Engineering, ICDE 2021, Chania, Greece, April 19-22, 2021*. IEEE, 2021, pp. 265–276. [Online]. Available: <https://doi.org/10.1109/ICDE51399.2021.00030>
- [37] A. Sinha, F. Fang, B. An, C. Kiekintveld, and M. Tambe, "Stackelberg security games: Looking beyond a decade of success," in *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, J. Lang, Ed. ijcai.org, 2018, pp. 5494–5501. [Online]. Available: <https://doi.org/10.24963/ijcai.2018/775>
- [38] M. Xiao, Y. Xu, J. Zhou, J. Wu, S. Zhang, and J. Zheng, "Aoi-aware incentive mechanism for mobile crowdsensing using stackelberg game," in *IEEE INFOCOM 2023 - IEEE Conference on Computer Communications, New York City, NY, USA, May 17-20, 2023*. IEEE, 2023, pp. 1–10. [Online]. Available: <https://doi.org/10.1109/INFOCOM53939.2023.10229079>
- [39] V. Conitzer and T. Sandholm, "Complexity results about nash equilibria," in *IJCAI-03, Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence, Acapulco, Mexico, August 9-15, 2003*, G. Gottlob and T. Walsh, Eds. Morgan Kaufmann, 2003, pp. 765–771. [Online]. Available: <http://ijcai.org/Proceedings/03/Papers/111.pdf>
- [40] C. Daskalakis, P. W. Goldberg, and C. H. Papadimitriou, "The complexity of computing a nash equilibrium," *SIAM Journal on Computing*, vol. 39, no. 1, pp. 195–259, 2009.
- [41] L. Xie, S. Meng, W. Yao, and X. Zhang, "Differential pricing strategies for bandwidth allocation with LFA resilience: A stackelberg game approach," *IEEE Trans. Inf. Forensics Secur.*, vol. 18, pp. 4899–4914, 2023. [Online]. Available: <https://doi.org/10.1109/TIFS.2023.3299181>
- [42] A. Sinha, P. Malo, and K. Deb, "A review on bilevel optimization: From classical to evolutionary approaches and applications," *IEEE Trans. Evol. Comput.*, vol. 22, no. 2, pp. 276–295, 2018. [Online]. Available: <https://doi.org/10.1109/TEVC.2017.2712906>
- [43] A. Marshall, *Principles of economics: unabridged eighth edition*. Cosimo, Inc., 2009.
- [44] T. Roughgarden, "Algorithmic game theory," *Communications of the ACM*, vol. 53, no. 7, pp. 78–86, 2010.
- [45] X. Ding, H. Wang, D. Zhang, J. Li, and H. Gao, "A fair data market system with data quality evaluation and repairing recommendation," in *Web Technologies and Applications: 17th Asia-Pacific Web Conference, APWeb 2015, Guangzhou, China, September 18-20, 2015, Proceedings 17*. Springer, 2015, pp. 855–858.
- [46] D. V. Winterfeldt and G. W. Fischer, "Multi-attribute utility theory: models and assessment procedures," *Utility, probability, and human decision making*, pp. 47–85, 1975.
- [47] L. R. Christensen, D. W. Jorgenson, and L. J. Lau, "Transcendental logarithmic utility functions," *The American Economic Review*, vol. 65, no. 3, pp. 367–383, 1975.
- [48] T. Başar and G. J. Olsder, *Dynamic noncooperative game theory*. SIAM, 1998.
- [49] J.-M. Lasry and P.-L. Lions, "Mean field games," *Japanese journal of mathematics*, vol. 2, no. 1, pp. 229–260, 2007.
- [50] S. Boyd, S. P. Boyd, and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [51] J. Castro, D. Gómez, and J. Tejada, "Polynomial calculation of the shapley value based on sampling," *Computers & OR*, vol. 36, no. 5, pp. 1726–1730, 2009. [Online]. Available: <https://doi.org/10.1016/j.cor.2008.04.004>
- [52] C. Dwork, "Differential privacy," in *Automata, Languages and Programming, 33rd International Colloquium, ICALP 2006, Venice, Italy, July 10-14, 2006, Proceedings, Part II*, ser. Lecture Notes in Computer Science, M. Bugliesi, B. Preneel, V. Sassone, and I. Wegener, Eds., vol. 4052. Springer, 2006, pp. 1–12. [Online]. Available: [https://doi.org/10.1007/11787006\\_1](https://doi.org/10.1007/11787006_1)
- [53] D. Dua and C. Graff, "UCI machine learning repository," 2017. [Online]. Available: <http://archive.ics.uci.edu/ml>