Enabling Secure Trustworthiness Assessment and Privacy Protection in Integrating Data for Trading Person-Specific Information

Rashid Hussain Khokhar[®], Farkhund Iqbal[®], Benjamin C. M. Fung[®], *Senior Member, IEEE*, and Jamal Bentahar[®], *Member, IEEE*

Abstract-With increasing adoption of cloud services in the e-market, collaboration between stakeholders is easier than ever. Consumer stakeholders demand data from various sources to analyze trends and improve customer services. Data-as-a-service enables data integration to serve the demands of data consumers. However, the data must be of good quality and trustful for accurate analysis and effective decision making. In addition, a data custodian or provider must conform to privacy policies to avoid potential penalties for privacy breaches. To address these challenges, we propose a twofold solution: 1) we present the first information entropy-based trust computation algorithm, IEB_Trust, that allows a semitrusted arbitrator to detect the covert behavior of a dishonest data provider and chooses the qualified providers for a data mashup and 2) we incorporate the Vickrey-Clarke-Groves (VCG) auction mechanism for the valuation of data providers' attributes into the data mashup process. Experiments on real-life data demonstrate the robustness of our approach in restricting dishonest providers from participation in the data mashup and improving the efficiency in comparison to provenance-based approaches. Furthermore, we derive the monetary shares for the chosen providers from their information utility and trust scores over the differentially private release of the integrated dataset under their joint privacy requirements.

Index Terms—Cloud computing, data mashup, data privacy, data trustworthiness, monetary valuation.

Rashid Hussain Khokhar and Jamal Bentahar are with the Concordia Institute for Information Systems Engineering, Concordia University, Montreal, QC H3G 1M8, Canada (e-mail: r_khokh@ciise.concordia.ca; bentahar@ ciise.concordia.ca).

Farkhund Iqbal is with the College of Technological Innovation, Zayed University, Abu Dhabi 144534, United Arab Emirates (e-mail: farkhund. iqbal@zu.ac.ae).

Benjamin C. M. Fung is with the School of Information Studies, McGill University, Montreal, QC H3A 0G4, Canada (e-mail: ben.fung@mcgill.ca).

Color versions of one or more of the figures in this article are available online at http://ieeexplore.ieee.org.

Digital Object Identifier 10.1109/TEM.2020.2974210

I. INTRODUCTION

ATA are the fuel of today's digital economy. Yet, data coming from a single source often fail to provide a complete picture for big data analytics. To answer complex queries, companies usually have to seek additional data from multiple sources. The emerging cloud paradigm data-as-a-service provides an ideal platform for data integration in order to serve data consumers' demands. However, business data often contain person-specific information. Mashing up personal data from different sources raises concerns on security, privacy, and data reliability. In the past decade, many trust models [6], [67] and frameworks [15], [57] have been proposed to evaluate and measure the security strength of cloud environments, but limited research considers the aspect of data reliability. In this article, we propose a cloud-based data integration solution that considers privacy protection, data trustworthiness, and fairness of profit distribution among data providers.

According to a recent survey [24], organizations in the U.S. estimate that 33% of their customer data are inaccurate. This skepticism about data elicits the increased risk of noncompliance and regulatory penalties. The study by IBM estimated that \$3.1 trillion of the U.S.'s GDP is lost due to poor quality data [64]. Organizations may mitigate these potential risks by taking appropriate measures regarding the quality of their data, leading to more reliable analysis and decision making. There is a line of research [13], [42] that focuses on exchanging data between multiple parties from the perspective of ensuring confidentiality and integrity. These works aim to provide prevention from unauthorized use and modification when data are in transit but do not verify data if any party provides false data. Our research perspective is to determine the trustfulness of private data held by dishonest data providers who may arbitrarily attempt to provide false data when trading person-specific information in the e-market for monetary benefits. Our proposed method can detect such behavior from dishonest data providers, who resemble adversaries under the covert security model [7]. In literature [3], [17], and [26], two protocols are discussed, namely, private set intersection (PSI) and PSI cardinality for privacy and data quality assessment. Freudiger et al. [27] claimed that these protocols are incurred from computational overhead and thus are not applicable to real-world scenarios. They proposed some protocols that operate on reduced dimensionality descriptions and

0018-9391 © 2020 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.

Manuscript received July 1, 2019; revised November 15, 2019 and February 4, 2020; accepted February 6, 2020. Date of publication March 2, 2020; date of current version November 13, 2020. This work was supported in part by the Research Cluster Award Fund R16083 and Research Incentive Funds R18055 and R19044 from Zayed University, in part by the Natural Sciences and Engineering Research Council of Canada under Discovery Grant RGPIN-2018-03872, and in part by the Canada Research Chairs Program under Grant 950-230623. Review of this manuscript was arranged by Department Editor P. Hung. (*Corresponding author: Benjamin C. M. Fung.*)

	Data Provider DP ₁			Data Provider DP_2			Data Provider DP ₃		
RecID	Age	Sex	Job	Sex	Education	Job	Age	Education	Job
1	39	М	Lawyer	M	Bachelors	Lawyer	45	Doctorate	Lawyer
2	50	М	Lawyer	М	Masters	Lawyer	50	Doctorate	Lawyer
3	38	М	Cleaner	M	12th	Technician	35	12th	Cleaner
4	53	М	Lawyer	M	Doctorate	Doctor	57	Masters	Lawyer
5	28	F	Cleaner	F	11th	Technician	28	11th	Cleaner
6	37	F	Welder	F	12th	Welder	37	11th	Welder
7	49	F	Painter	F	12th	Cleaner	49	12th	Painter
8	59	М	Doctor	F	Doctorate	Doctor	66	Doctorate	Doctor
9	31	F	Painter	M	12th	Welder	27	12th	Painter
10	42	М	Technician	М	Bachelors	Technician	42	Bachelors	Technician
11	37	М	Lawyer	M	Masters	Lawyer	38	Masters	Lawyer
12	30	М	Lawyer	М	Masters	Lawyer	28	Bachelors	Lawyer

TABLE I RAW DATA OWNED BY THREE DATA PROVIDERS



Fig. 1. Taxonomy trees.

so can be scalable to large datasets. It is a challenging problem to evaluate the trustfulness of private data held by untrusted data providers. In this article, we study the problem of untrusted data providers holding overlapping attributes on a person-specific dataset. We illustrate the problem in the following example.

Example 1: Suppose that there is a cloud-based data market, where data consumers can place their data mining requests and data providers compete with each other to contribute their data with the goal of fulfilling the requests for monetary reward. Consider the 12 raw data records in Table I, where each record corresponds to the personal information of an individual. The three data providers own different yet overlapping sets of attributes over the 12 records.

Since the data providers collect data from different channels, it is quite possible that their data conflict with each other, as illustrated in Table I. According to the predefined generalization hierarchy of the attributes in Fig. 1, the individuals in the table can be generalized to two groups: *Non-Technical* and *Technical*. Suppose that a data consumer wants to perform a data analysis that depends on the *Non-Technical* and *Technical* groups. Yet, the inconsistent, conflicting, or even inaccurate data may mislead the analysis result. For example, DP₁ and DP₃ state that the individuals in {Rec#3,5} are *Cleaner*, while DP₂ states that they are *Technician*. A similar conflict can be seen in Rec#9, where DP₁ and DP₃ provide the *Job* as *Painter*, and DP₂ provides the *Job* as *Welder*. In this example, the *Job* attribute on {Rec#3,5,9} has two different values that are categorized as *Non-Technical* and *Technical*, respectively. These inconsistencies significantly impact the quality of data analysis.

Presumably, the data providers would have missing values on some attributes, although the same set of records is identified by executing the secure set intersection protocol [3] on the globally unique identifiers [53], [54]. Instead of avoiding participating in the data mashup process, they would prefer to impute missing values by using the machine learning methods appropriate for their datasets. The properties of a dataset such as low-dimensional or high-dimensional data, single-type or mixed-type data, or linearly separable or nonlinearly separable data are a crucial factor before choosing the imputation method. The data providers' decision whether to use a single imputation method or multiple imputation methods is conditional on their missing data. We evaluate the robustness of our approach when an acquisitive data provider employs a machine learning method for imputation of missing data.

In the context of quantifying monetary value through sharing person-specific data, the data providers first must do the valuation of personal data, but there is no determined market price [56], [62] for person-specific data that can be taken as a proxy for the valuation. It is also well acknowledged from the existing literature [25], [58] that there is no commonly agreed methodology for valuing personal data. However, in the e-market, many companies actively collect personal information by providing monetary rewards to their customers. In this article, we incorporate the Vickrey-Clarke-Groves (VCG) auction mechanism for the valuation of data providers' attributes. We reason that it is a dominant strategy, where no data provider has an incentive to lie about his true valuations. In addition, private data often encode privacy-sensitive information related to individuals that need to be protected when integrating data from the competing data providers. In this article, we adopt differential privacy [22] because it provides strong privacy guarantees to an individual independently of an adversary's background knowledge, in contrast to underlying assumptions in syntactic privacy models [47], [51], [66] about an adversary's knowledge.

Contributions: We propose a novel solution to address the critical issues of data trustworthiness, privacy protection, and profit distribution for cloud-based data integration services. The data trustworthiness problem has been studied in [49], [50], and [69] applications of sensor networks. The provenance-based approach has been used in [16] and [50] to evaluate the trust-worthiness of network nodes and data items. This approach is primarily used to collect evidence about where the data originate and how the data generate. In this article, we are not concerned about the high degree of the instrumentation of customers' private data, which are collected by data providers. However, our proposed approach makes novel use of information entropy to verify the correctness of data from untrusted data providers and also to preserve the privacy of customers' data held by data providers when evaluating the trustworthiness of the providers.

PSI-based approaches allow multiple parties to jointly compute the intersection of their private data without revealing any additional information to either side [75]. These approaches are suitable for privacy-preserving distributed data mining, in which multiple data custodians compute a function based on their inputs without sharing their data with others. In this article, we focus on privacy-preserving data publishing (*PPDP*) in a distributed setting, where the data providers wish to integrate their data for better information utility. However, the data integration necessitates that under the specified privacy constraints, no data provider should learn any additional information other than necessary information. We summarize our contributions as follows.

- Our proposed method, *IEB_Trust*, is the first entropybased trust computation method that enables secure trustworthiness assessment and incorporates fairness in the verification process to restrict dishonest data providers from participation in the next phase for integrating data.
- We compare our proposed method with a closely related method. Results suggest that our entropy-based trust computation algorithm is capable of significantly improving runtime efficiency.
- 3) We evaluate the robustness of our method when an acquisitive data provider adopts machine learning techniques to substitute missing values on their own data and claim them as original data collected from customers to compete with the other participating data providers.
- 4) We define the procedure for setting the price on personspecific attributes in trading personal information from data providers based on the VCG mechanism.
- 5) We integrate data from chosen data providers using differentially private anonymization based on generalization (*DistDiffGen*) [53] and analyze the impacts of privacy protections and trust scores on data providers' monetary value.

The rest of this article is organized as follows. In Section II, we provide an overview of the trust mechanism and the problem statement. In Section III, we review the related work. In Section IV, we discuss the trust aspects, imputation methods, and privacy models. In Section V, we present our proposed solution. In Section VI, we compare our proposed method and provide empirical study to analyze the trustworthiness of each data provider and further analyze its impact along with the ϵ -differential privacy protection on a data provider's monetary value. Finally, Section VII concludes this article.

II. TRUST MECHANISM

In this section, we first provide an overview of our trust mechanism and then formally define the research problem.

A. Overview of the Trust Mechanism

Fig. 2 provides an overview of our trust mechanism, in which data providers, data consumers, and cloud service providers (CSPs) are the main entities. Data providers collect person-specific information from customers and intend to participate in the data mashup for generating more profit by competing with peer data providers, data consumers perform data analysis on the received data, and the CSP is a semitrusted arbitrator between data providers and data consumers. The CSP manages three key services: authentication, mashup coordination, and data verification. These services are run on a cloud server (CS) by the CSP. First, each data provider has to pass the authentication



Fig. 2. Trust mechanism.

phase to prove their identity. Second, data consumers submit their data requests to the CSP. In this article, we assume that a data consumer runs a classification analysis on its requested attributes by a supervised machine learning method. A resource queue is built by the mashup service to manage data requests from a data consumer, which is accessible only to authenticated data providers. Third, data providers register their available data attributes on the registry hosted by the mashup service; each data attribute is assigned a sequence number based on its arrival. Fourth, the verification process is run to detect false or incorrect data and to determine the trustworthiness of each data provider. Fifth, this process results in determining the accepted data providers. Sixth, the CSP connects the group of accepted data providers with the data consumer to serve its demand. This is done by the mashup service that determines the group of data providers, whose data can collectively fulfill the demand of a data consumer. Seventh, the data providers quantify their costs and benefits using joint privacy requirements and integrate their data over the cloud. Finally, the anonymous integrated data are released to the data consumer.

B. Problem Statement

We describe our problem as follows. There are three main entities discussed in our trust mechanism: data providers, data consumers, and a CSP. *Data verification* service runs on a CS, which is managed by the CSP. The purpose of this service is to verify the correctness of data. The CSP is a semitrusted arbitrator who would not have access to customers' private data, which is held by the data providers. Data providers are considered to be dishonest, meaning that they may arbitrarily attempt to provide false data because they are acquisitive in competing with others in the e-market. The behavior of such data providers is similar to adversaries in the covert security model.

Suppose that data providers $DP_1,..., DP_n$ own private data tables D_1, \ldots, D_n , respectively. Each record in the data table belongs to a unique individual. All explicit identifiers of an individual, such as name, social security number (SSN), and account number, have been removed. Each D_i is defined over a set of attributes $\mathcal{PA}_i = \{A_1, \ldots, A_d\}$. We assume that the data providers hold overlapping attributes for the same set of records identified by executing the secure set intersection protocol [3], [54] on the globally unique identifiers RecID. We require $\forall \mathcal{PA}_i \exists \mathcal{PA}_j$ such that $\mathcal{PA}_i \cap \mathcal{PA}_j \neq \emptyset$, where $i \neq j$, and $\mathcal{PA} = \{\mathcal{PA}_1, \dots, \mathcal{PA}_n\}$. In addition, each D_i contains an A^{cls} attribute for classification analysis, which is shared among all the data providers. Each $A_{\mathcal{I}}$ is either a categorical or a numerical attribute, but A^{cls} is required to be categorical. A data consumer submits a data request $\operatorname{Req} A = \{\operatorname{Req} A_1, \dots, \operatorname{Req} A_m\}$ for classification analysis. We assume that each data provider has $\mathcal{PA}_i \subseteq \operatorname{Req} A$ to serve the demand of a data consumer. The goal of this trust computation is to restrict dishonest data providers from participation in the data mashup process when their trust scores drop below a certain threshold.

Problem 1 (Trust Computation): Given multiple personspecific raw data tables D_1, \ldots, D_n from data providers DP_1, \ldots, DP_n and a set of requested attributes $ReqA = \{ReqA_1, \ldots, ReqA_m\}$ for classification analysis from a data consumer, the research problem is to verify the correctness of data on the submissions of the overlapping set of attributes $\mathcal{PA}_i = \{A_1, \ldots, A_d\}$ on the same set of records from each data provider DP_i , where $\mathcal{PA}_i \cap \mathcal{PA}_j \neq \emptyset \forall \mathcal{PA}_i \exists \mathcal{PA}_j$ and $i \neq j$ and to compute the trust score TS_{DP_i} of each data provider.

In the context of data privacy, the data providers want to integrate their data in a way such that no data provider should learn any additional information about the others as a result of data integration. After the completion of trust computation, the data providers DP_1, \ldots, DP_n attain a mutually exclusive set of attributes $\mathcal{PA}_i = \{A_1, \ldots, A_d\}$ over the same set of records for

data integration. That is, $\mathcal{PA}_i \cap \mathcal{PA}_j = \emptyset$ for any $1 \leq i, j \leq n$. We assume that for each attribute $A_{\mathcal{J}} \in \mathcal{P}A_i$, a taxonomy tree is provided that defines the hierarchy of values in $\Omega(A_{\mathcal{J}})$, where $\Omega(A_{\mathcal{J}})$ represents the domain of $A_{\mathcal{J}}$. Data providers require doing their attributes' valuations for price setting and jointly setting up the privacy requirements, such as privacy budget ϵ and specialization level h for a ϵ -differential privacy model, before data integration. They wish to derive their monetary shares from the information utility of anonymous integrated data \hat{D} for classification analysis and their trust scores.

Problem 2 (Monetary Share Under ϵ -Differential Privacy *Mechanism*): Given multiple raw data tables D_1, \ldots, D_n containing mutually exclusive sets of attributes $\mathcal{PA}_i =$ $\{A_1, \ldots, A_d\}$, i.e., $\mathcal{PA}_i \cap \mathcal{PA}_j = \emptyset$ for any $1 \le i, j \le n$ over the same set of records, and a data request ReqA = $\{\operatorname{Req}A_1,\ldots,\operatorname{Req}A_m\}$ from a data consumer for classification analysis, the research problem is to derive the monetary share of each DP_i from their information utility and trust scores over the differentially private release of integrated dataset D under the joint privacy requirements and attributes' valuations.

Several companies, such as Acxiom, AnalyticsIQ, Dataline, and Expedia, collect user data, including demographic, financial, retail, social, and travel information from multiple sources with the goal of serving different market needs [1]. Our research problem can be generalized to other similar companies who face trustworthy or quality data issues [24] and whose business models are primarily based on sharing person-specific information.

III. RELATED WORK

In this section, we summarize the literature of the following related areas: data trustworthiness and auction-based pricing, cryptographic primitives, and differentially private anonymization techniques.

A. Data Trustworthiness and Auction-Based Pricing

Different trust models, frameworks, and techniques have been proposed to address the problem of data trustworthiness. Bertino and Lim [11] proposed a framework that consists of two key components. The first component is based on the concept of data provenance, in which information relies on the origin of data for computation of trust scores. The second component undertakes the notion of confidence policy, in which query results are filtered based on the specified confidence range for use in certain tasks. Dai et al. [16] proposed a provenance-based model, in which they evaluated the trustworthiness of data items based on the aspects of data similarity, path similarity, data conflict, and data deduction. Benjelloun et al. [8] introduced databases with uncertainty and lineage, in which they combined the concept of *lineage* and *uncertainty* for querying in probabilistic databases.

There are studies related to data trustworthiness in missioncritical applications [49], [69]. Tang et al. [69] proposed trustworthiness analysis for sensor networks in cyber-physical systems to eliminate false alarms that occur due to random noise or defective sensors. They validated events by using a graph-based filtering approach. However, their method does not deal with coordinated attacks, where a fraction of sensing nodes are compromised by malicious attackers. Lim et al. [49] addressed this challenge by adopting a game-theoretic approach based on the Stackelberg competition for defending the network against false data injection. They assessed trust scores for both data items and network nodes using the cyclic framework proposed in [50]. This framework is based on the interdependence property between data items and their associated network nodes in which trust scores are computed using two types of similarity functions. First, *value similarity* is derived from the principle that the more that similar values refer to the same event, the higher the trust scores. Second, provenance similarity is based on the principle that the more that different data sources are with similar data values, the higher the trust scores. Mainly, the approaches presented in the above works fall under the category of workflow provenance. In contrast, we are not concerned about the higher level of instrumentation at the data collection phase by data providers because it is not practically efficient to determine the data provenance in the e-market. Furthermore, the above works mainly focus on similarity functions for trust computation but do not consider privacy protection for data trustworthiness. We propose an approach that makes novel use of information entropy to verify the correctness of data in a multiple data providers' scenario, where a semitrusted arbitrator cannot derive any customers' private data when evaluating the trustworthiness of the participating data providers.

Karabati and Yalcin [41] studied the challenge of pricing with short-term capacity allocation decisions for multiple products in a single-supplier multiple-buyer scenario. They proposed an iterative auction mechanism with monotonically increasing prices to maximize the profit of a supplier. Li et al. [48] presented dynamic pricing strategies for resources allocations in cloud workflow systems. Their proposed reverse-auction-based mechanism allows resource providers to change the prices during the auction, depending upon their trading situation, to improve the efficiency of resource utilization as well as the competitiveness. Wu et al. [72] employed a VCG auction to implement a dynamic pricing scheme for multigranularity service composition. They considered both coarse-grained and fine-grained services for composition. In their approach, service providers bid for services of different granularities in the composite service, whereas a recipient of the bids decides a composition that minimizes the overall cost while satisfying quality constraints. They solved the problem of winner determination by an integer programming model. In this article, we define the procedure for the valuations of data providers attributes based on the VCG mechanism.

B. Cryptographic Primitives

PSI is a cryptographic primitive that was first formally defined in [26]. The protocols for PSI allow two parties, holding sets Aand B, to compute the private intersection without revealing to each other any additional information from their respective sets. At the end of the protocol, either one or both parties may learn the size of the intersection, depending on the application. Since its inception, many variants have been proposed in an attempt to speed up PSI computation, including garbled Bloom filters [20], [33], server-aided computations [19], [39], [40], and computational optimizations [46], [59], [61].

Oblivious transfer (OT) is one of the fundamental primitives in cryptography and has been extensively used for secure multiparty computation. Particularly, the most efficient OTs were introduced by Pinkas *et al.* [61] and further strengthened in [46], [59], and [60]. Kolesnikov *et al.* [46] proposed a batched related-key oblivious pseudorandom function (OPRF) protocol to improve the performance of semihonest secure PSI. They achieved a 1-out-of-n OT of random messages for an arbitrarily large n at nearly the same cost as 1-out-of-2 in [35]. The new OPRF construction of Pinkas *et al.* [59] is similar to Kolesnikov *et al.* [46] except in handling error correcting code. Kolesnikov *et al.* [46] demonstrated that their protocol outperforms Pinkas *et al.* [60] in almost as many settings, particularly for the long bit length of input and large values of the input size.

In practice, the OT-based protocols are much faster than the random garbled Bloom filter-based protocols for larger set sizes, yet these protocols do not have the lowest communication cost [46]. One desirable property is to achieve the fairness that ensures either all the parties of a group learn the output of the computation or none do [39]. This is not the case with standard approaches to PSI. Our solution to the problem is different from several PSI-based approaches, in which the intention is to achieve both privacy and security simultaneously. These approaches are suitable for different motivating applications in private data mining, online recommendation services, and genomic computations. In our approach, we maintain confidentiality and integrity by exchanging only an encrypted information gain message and its keyed hash between a data provider and the CS, based on a random challenge (i.e., attribute request) of the CS, instead of exchanging encrypted individual data items. This apparently reduces the overhead of communication. In addition, we do not rely on the server to perform the computation on clients' private data. In the context of privacy, PSI protocols enable parties to privately know the result from their intersection, but the total information is not published for data analysis [75]. However, we intend to securely integrate person-specific data from multiple data providers and to release differentially private data for classification analysis.

C. Differentially Private Anonymization Techniques

Differential privacy is increasingly being accepted as the cornerstone of privacy protection by domain experts due to its robustness and rigorous mathematical definition. In the literature, two settings, namely *interactive* and *noninteractive*, are mainly discussed regarding utilization of the privacy budget ϵ . The primary difference is that in the interactive setting [22], [28], [73], [74], the data custodian holds the raw data, and a data analyst poses a set of queries in real time, for which the data custodian provides differentially private answers. Each query would utilize a fraction of ϵ incrementally to produce a noisy answer. When the entire privacy budget has been depleted, a data analyst would not be able to get the answer by querying the database. On the other hand, in the noninteractive setting, the data custodian first anonymizes its raw data by utilizing

the entire privacy budget. Later, the anonymous (ϵ -differentially private) data releases to the data analyst, who would perform an analysis without any constraints on the data usage. This approach is widely known as PPDP [30], which is more appropriate in many real-life data sharing scenarios because of the flexibility for a data analyst to perform an analysis without back and forth querying of the database. In this article, we focus on the noninteractive setting for a differentially private release of data in a distributed setup.

The group of works [4], [53] based on distributed approaches are suitable for multiple parties, whose prime concern is to integrate their data in a way that no party could learn any additional information about the other party as a result of data integration. Mohammed et al. [53] proposed an algorithm, called DistDiffGen, in which data are vertically partitioned among multiple parties in a distributed setup. It allows two parties to securely integrate their person-specific data while maintaining necessary information to support data utility. Each party in this setup owns a mutually exclusive set of attributes over the same set of records. A similar problem has also been studied by Alhadidi et al. [4], where data are horizontally partitioned among two parties. Each party in this setup owns a disjoint set of records over the same set of attributes. In this article, we employ DistDiffGen [53] for a distributed setup with an extension for multiple data providers to achieve ϵ -differential privacy. There are existing works that allow data integration for horizontally partitioned databases [37], [55] and vertically partitioned databases [29], [36], [54] under the privacy constraints in a distributed setup. These works are based on syntactic privacy models, which are vulnerable to certain attacks such as minimality attack [71], composition attack [32], and deFinetti attack [44]. Therefore, we adopt differential privacy [22] because it provides strong privacy guarantees against such attacks. Whereas existing work [43] proposed a privacy-preserving data mashup model that allows the collaboration of multiple data providers for integrating their data and derives the contribution of each data provider by evaluating the incorporated cost factors, in our article, we derive the monetary shares for the chosen data providers from their contribution to information utility over the differentially private integrated data for classification analysis and their trust scores.

IV. PRELIMINARIES

In this section, we first present the principles that are crucial for establishing trust. Next, we discuss methods for imputation of missing data, and finally, we discuss privacy models.

A. Trust Aspects

Trust is a critical aspect of decision making in e-commerce. Trust principles are a part of many service-oriented-architecturebased models, where participants in the system want to do interactions for service delivery and use [2]. We review the principles that are crucial for trust establishment. First, entities should be identified [38] as they have claimed. In the world of the Internet, where entities are physically isolated, they may have real identifies or may use fake identities to show their existences in their interactions. Authentication is a way of validating entities by the use of usernames and passwords, tokens, or digital certificates before granting them access to the resources or applications [12]. Second, it is crucial for trust formation to initialize new entities with trust rates. This process is called trust bootstrapping. Third, when one entity trusts another entity's decision, there is a risk of an undesirable outcome due to some degree of uncertainty and dependence [45]. The risk is considered to be a prerequisite before trusting the trustee's behavior. The entities who are involved in an interaction should comply with the norms and rules of trust to avoid penalties for violation. Fourth, trust rates are of two types: local and global [70]. Local trust rating refers to a personalized score, in which each trustee would have different scores from the trustors. Global trust rating provides a unique score about the entity (trustee) independently of who are the entities (trustors) participating in the evaluation. Global trust rating often requires the trusted third party services to collect feedback from the trustors about trustees and compute the trust rates. Last, security and privacy are the main components for trust establishment. Trust is required when there is uncertainty; it has widely been accepted that perfect security does not exist, even though security measures are necessary to gain trust in many circumstances [10]. Customers who place their orders online and submit private information in the form of their name, address, and credit details necessitate that their private information should not be disclosed or shared by any means with untrusted parties. Building a trust relationship requires protection of customers' privacy in online transactions. We pay attention to some of the aforementioned principles for establishing trust on the data providers in the context of our trust mechanism.

B. Methods for Imputation of Missing Data

There are different types of missing data [34], such as missing at random (MAR), missing completely at random (MCAR), and missing not at random (MNAR). MAR refers to the probability of missing data of an attribute on other present observations of attributes in the dataset, but not on the attribute's own value, whereas MCAR occurs when there is no dependence on the attribute value itself or any other attribute in the dataset. And the special case MNAR occurs when the missing data meet neither the condition defined in MAR nor MCAR. In this special case, missing values in MNAR cannot be imputed by using other present observations of attributes.

There is extensive research [5], [9], [76], [77] done on machine learning methods such as hot-deck imputation, mean imputation, regression imputation, k-nearest neighbor (kNN) imputation, and random forest imputation. Hot-deck imputation is a technique for replacing missing values of a nonrespondent on one or more attributes with the most similar characteristics to a respondent [5]. This method has been used in practice, but the theory is not as well developed. Mean imputation is a technique used for replacing missing values of a numerical attribute by the average value, and for a categorical attribute by the mode, i.e., most frequent value. This method is quite simple, but it is not suitable for multivariate analysis. Regression imputation first builds a model from the observed data; then, predictions for the incomplete cases are calculated under the fitted model to replace the missing data [77]. The drawback of the regression model is that all predicted values fall directly on the regression line, which decreases variability. Random forest is a type of ensemble learning method [76]. It is used widely for classification and regression tasks. The learning process of a random forest algorithm is based upon the bootstrap aggregation technique, in which a specified number of trees are trained on a given dataset. As the random forest is built upon multiple decision trees, intrinsically, it uses the same approach for attribute selection measures such as information gain, gini index, and gain ratio of decision trees. Random forest can deal with missing values with different types of variables. kNN imputation is an efficient approach for replacing missing values on some records by computing another value from similar examples in the given dataset [9]. kNN computes the similarity by using a distance metric, such as Euclidean distance. k is a positive integer, when k = 1, the object is simply assigned to the class of that single nearest neighbor. When k > 1, the object is assigned to the class that appears most frequently within the k-subset. kNN generally produces good quality predictions, but the computation cost is high because of computing distances.

C. Privacy Models

In the literature, there are two types of models apprehended: *syntactic* and *semantic*. Syntactic models, such as *K*-anonymity [66], protect against identity disclosure, *l*-diversity [51] protects from homogeneity attacks, and *t*-closeness [47] is an extension of *l*-diversity, in which the distribution of sensitive attribute values for privacy protection is further refined. Differential privacy [22] is a semantic model that is more robust against the aforementioned attacks. It provides strong privacy guarantees to an individual independently of an adversary's background knowledge. The intuition of differential privacy is that individual information is not revealed from the output of the analysis in the anonymized data. In other words, it is insensitive whether an individual record is present in the input dataset or not. It is mathematically defined as follows.

Definition IV.1 (ϵ -differential privacy [22]): A sanitization mechanism M provides ϵ -differential privacy, if for any neighboring datasets D and D' differing by at most one record (i.e., symmetric difference $|D \triangle D'| \leq 1$), and for any possible sanitized dataset \hat{D} , we have

$$\Pr[M(D) = \hat{D}] \le e^{\epsilon} \times \Pr[M(D') = \hat{D}]$$

where the probability is taken over the randomness of the M.

 ϵ is the privacy budget that is specified by the data custodian. A smaller value of ϵ results in stronger privacy protection but produces lower data utility. Conversely, a larger value of ϵ results in weaker privacy protection but yields higher data utility.

The Laplace mechanism and the exponential mechanismare the canonical examples of a differentially private mechanism. A standard mechanism to achieve differential privacy is to add random noise to the outcome of the analysis for providing privacy protection. The calibration of noise is done according to the sensitivity of the function f.

Definition IV.2 (Sensitivity) For any function $f: D \rightarrow \mathbb{R}^d$, the sensitivity of f is

$$\Delta f = \max_{D,D'} ||f(D) - f(D')||_1 \tag{1}$$

for all D, D' differing at most by one record.

The sensitivity of a function does not depend on the data but instead produces an upper bound to how much noise we must add to the true output to preserve privacy. Suppose that function f answers count queries over a dataset D. Then, Δf is 1 because f(D) can differ at most by 1, due to the addition or removal of a single record.

Laplace mechanism: Dwork et al. [22] proposed the Laplace mechanism. It is appropriate when the output of function f is a real value, and f should perturb its output with a noisy answer to preserve privacy. The noise is calibrated based on the privacy parameter ϵ and the sensitivity of the utility function Δf . Formally, the Laplace mechanism takes as inputs a dataset D, the privacy parameter ϵ , and a function f and outputs $\hat{f(D)} = f(D) + \text{Lap}(\lambda)$, where $\text{Lap}(\lambda)$ is a noise drawn from the Laplace distribution with the probability density function $\Pr(x|\lambda) = \frac{1}{2\lambda} \exp(-|x|/\lambda)$. The variance of this distribution is $2\lambda^2$, and the mean is 0.

Exponential mechanism: McSherry and Talwar [52] proposed the exponential mechanism. It is appropriate for situations in which it is desirable to choose the best response, because adding noise directly to the count can eradicate its value. Given an arbitrary range \mathcal{T} , the exponential mechanism is defined with respect to a utility function $u: (D \times \mathcal{T}) \to \mathbb{R}$ that assigns a real-valued score to every output $t \in \mathcal{T}$, where a higher score means better utility. The exponential mechanism induces a probability distribution over the range \mathcal{T} and then samples an output t. Suppose $\Delta u = \max_{\forall t, D, D'} |u(D, t) - u(D, t)|$ to be the sensitivity of the utility function. The probability associated with each output t is proportional to $\exp(\frac{\epsilon u(D,t)}{2\Delta u})$.

V. PROPOSED SOLUTION

In this section, we provide a solution to address the concerns of stakeholders on data trustworthiness, privacy protection, and profit distribution in the online market for trading person-specific data. Section V-A presents our proposed *IEB_Trust*, an information entropy-based trust computation algorithm to restrict dishonest data providers from participation in the data mashup process and to assess the trustworthiness of each data provider. Section V-B discusses security properties. Section V-C provides an analysis of the IEB_Trust algorithm. Section V-D provides an evaluation of learner models. Section V-E provides an auction mechanism for price setting among data providers who own multiple attributes. Section V-F presents an algorithm for privacy protection by which data providers can determine the impact of anonymization on data utility for classification analysis. Section V-G discusses how the chosen data providers can quantify their monetary value.

A. Trust Computation

In Section II-B, we state the problem where the challenge is to verify the correctness of data from untrusted multiple data providers, who own overlapping attributes for the same set of records. We assume that the data providers are competitors, who intend to maximize their profits. The data providers consider as dishonest anyone who may arbitrarily attempt to provide false data to get a larger monetary share from their participation. To address this problem, we propose a novel algorithm that adopts information entropy for secure trustworthiness assessment of acquisitive data providers. Information entropy has been widely used in machine learning tools and decision-making systems. Compared to the existing work on data trustworthiness [49], [50], [69], our proposed algorithm not only detects false or incorrect data from a dishonest data provider during the verification process, but also preserves the privacy of customers' data owned by a data provider. Furthermore, our method provides better runtime efficiency over provenance-based approaches [16], [50].

Algorithm 1 presents our approach in more detail. A CSP runs this algorithm on a CS. Consider multiple data providers DP_1, \ldots, DP_n , who own private data tables D_1, \ldots, D_n having overlapping attributes for the same set of records identified by the common record identifier RecID [3], [54]. First, the CS and each DP_i mutually authenticate each other and derive ks_i symmetric keys for all $i \in I$ by the mutual authentication protocol [18] for the secure exchange of messages. Each DP_i has its own ks_i to answer the CS's queries. Second, a data consumer submits a data request $\operatorname{Req} A = \{\operatorname{Req} A_1, \ldots, \operatorname{Req} A_m\}$ to the CS. Third, each data provider DP_i submits an available set of attributes $\mathcal{PA}_i =$ $\{A_1, \ldots, A_d\}$, where $\mathcal{P}\mathcal{A}_i \subseteq \operatorname{Req}A$, to the CS. We assume that, initially, all the participating data providers have "zero" in their trust scores (Line 3). ϵ' is the allocated privacy budget to consume for each requested attribute. A resource queue is created by the mashup service for m requested attributes, where each attribute $A_{\mathcal{J}} \in \mathcal{P}\mathcal{A}_i$ of a corresponding data provider is registered with its arrival sequence (Line 9).

Fourth, the verification process is run to determine the trustworthiness of each data provider. In the first round, the CS successively selects one attribute $\operatorname{Req} A_x'$ uniformly at random without replacement over a domain of m requested attributes and sends an encrypted challenge $E(ks_i, \text{Req}A_x')$ to the corresponding data providers DP_1, \ldots, DP_n , who own common attribute $A_{\mathcal{J}}$. Prior to responding to this challenge, each DP_i decrypts to retrieve $\operatorname{Req} A_x'$ computes information gain on the challenge attribute in Line 16, denoted by $\mathcal{G}_{A_{\mathcal{I}}}^{(1)}$ (refer to Section V-A1 for details), according to (4) [63], and then adds noise to a true output. Then, DP_i encrypts the message $\psi^{(1)} \leftarrow E(ks_i, \mathcal{G}_{A_{\mathcal{I}}}^{(1)'})$ and computes tags $\Upsilon^{(1)} \leftarrow \mathcal{S}(k_h, \psi^{(1)})$ by using keyed hash-based message authentication code (HMAC) in Line 17. The CS receives the concatenated message, tag, and identity $\psi^{(1)} \| \Upsilon^{(1)} \| DP_i$ on his challenge from each data provider. Then, the CS computes the comparison to determine the majority candidates by invoking procedure findMajCand($\psi^{(1)} \| \Upsilon^{(1)}$, size) in Line 19, where *size* indicates the number of data providers who own the requested attribute. This procedure returns majority candidate Maj $_{Cand}^{R(1)}$. In the second round, the CS generates \mathcal{K} random

IDs for the requested challenge $\operatorname{Req} A_x'$, i.e., picked in the first round, from $|D_i|$ records, and then generates \mathcal{P} pairs of values for $\operatorname{Req} A_x'$ and A^{cls} attributes. The CS sends another challenge to each DP_i by concatenating the encrypted \mathcal{K} random IDs and \mathcal{P} pairs of values as $E(ks_i, \mathcal{K}, \text{Reg}A_x') || E(ks_i, v_{x'}, v_{cls})$. DP_i decrypts to retrieve \mathcal{K} record IDs and \mathcal{P} pairs of values. DP_i concatenates ${\mathcal K}$ records and ${\mathcal P}$ pairs of values received from the CS. DP_i computes $\mathcal{G}_{A_{\mathcal{T}}}^{(2)}$ on the concatenated version and then adds noise to a true output, encrypts it as $\psi^{(2)} \leftarrow E(ks_i, \mathcal{G}_{A_{\mathcal{T}}}^{(2)'})$, and computes the tag as $\Upsilon^{(2)} \leftarrow \mathcal{S}(k_h, \psi^{(2)})$. The CS receives $\psi^{(2)} \| \Upsilon^{(2)} \| DP_i$ on the second round challenge from the corresponding data providers in Line 28. The CS again invokes procedure findMajCand($\psi^{(2)} \| \Upsilon^{(2)}$, size) to determine the majority candidates in Line 29. This process repeats α times. In Line 33, an intersection of both the rounds is computed to determine Maj_{Cand}.

Candidates whose scores match on the majority are considered as Qualified, denoted by Qual_{DP}, who gain a positive weight γ in their trust scores TS_{DP_i} . Alternatively, candidates whose scores do not match are considered as Nonqualified, denoted by $UnQual_{DP_i}$. Subsequently, $UnQual_{DP_i}$ is penalized with a negative weight $-\gamma$ in their trust scores TS_{DP_i} . When only a single data provider responds to the CS challenge of $\operatorname{Req} A_x'$, it is accepted based on his existing trust score $TS_{DP_a} > 0$. However, in this case, the trust score does not increase for that data provider. When a data consumer request for an attribute, which is not fulfilled by the participating data providers, then that attribute is excluded from the verification process, and the data providers gain no monetary value from it. The comparison is performed (Line 45) to select one candidate (or data provider) on each attribute from the qualified data providers Qual_{DP_n} based on their arrival sequences (using first-come first-served rule). If the final aggregated trust score of any data provider becomes < 0, that data provider drops from the final selection for the data mashup, and the attributes initially belonging to him are subsequently reassigned to other qualified data providers that appear next in the arrival sequences. The algorithm terminates when there is no more attribute for verification.

1) Computation of Information Gain: We use information gain as a criterion for splitting attributes [63] based on the concept introduced by Claude Shannon on information theory [68]. We compute information gain on an individual attribute $A_{\mathcal{J}} \in \mathcal{P}A_i$ of each data provider in the presence of a shared class attribute A^{cls} on raw data. Let $D^{\tau} \subseteq D_i$ denote a subset of the data table D_i . Suppose that the attribute A^{cls} has \mathcal{C} distinct values. Let $A_{i,D^{\tau}}^{\text{cls}}$ be the set of records of class A_i^{cls} in D^{τ} . Let $|D^{\tau}|$ and $|A_{i,D^{\tau}}^{\text{cls}}|$ denote the number of records in D^{τ} and $A_{i,D^{\tau}}^{\text{cls}}$, respectively. The entropy on the data table D^{τ} is computed as follows:

$$E(D^{\tau}) = -\sum_{i=1}^{\mathcal{C}} \operatorname{Pr}_i \times \log_2 \operatorname{Pr}_i$$
(2)

where \Pr_i is the probability that an arbitrary record in D^{τ} belongs to class A_i^{cls} . It is estimated by $\frac{|A_{i,D^{\tau}}^{\text{cls}}|}{|D^{\tau}|}$.

We can further partition the records in D^{τ} on the attribute $A_{\mathcal{J}}$. If $A_{\mathcal{J}}$ is discrete valued, then one branch is grown for each known value of $A_{\mathcal{J}}$. On the other side, if $A_{\mathcal{J}}$ is continuous valued, then two branches are grown, corresponding to $A_{\mathcal{J}} \leq$ splitpoint and $A_{\mathcal{J}} >$ splitpoint. It is calculated by the following equation:

$$E_{A_{\mathcal{J}}}(D^{\tau}) = \sum_{j=1}^{\mathcal{V}} \frac{|D_j^{\tau}|}{|D^{\tau}|} \times E(D_j^{\tau}).$$
(3)

Finally, we can compute the information gain $\mathcal{G}_{A_{\mathcal{J}}}$ on the chosen attribute $A_{\mathcal{J}}$ of each data provider DP_i as follows:

$$\mathcal{G}_{A_{\mathcal{T}}} = E(D^{\tau}) - E_{A_{\mathcal{T}}}(D^{\tau}). \tag{4}$$

2) Differentially Private $\mathcal{G}_{A_{\mathcal{J}}}$: Given a privacy budget ϵ' , the sensitivity of the utility function (Δf) is 1, and a true computed $\mathcal{G}_{A_{\mathcal{J}}}$. We add independently generated noise from the Laplace distribution $\text{Lap}(1/\epsilon')$ to a true computed $\mathcal{G}_{A_{\mathcal{J}}}$ to have a differentially private version of (4):

$$\mathcal{G}'_{A_{\mathcal{J}}} = \mathcal{G}_{A_{\mathcal{J}}} + \operatorname{Lap}(1/\epsilon').$$
 (5)

3) Discretization: We use equal-width method to discretize a continuous-valued attribute $A_{\mathcal{J}}$ into K intervals of equal size. The min_{val} and max_{val} parameters are used for defining the boundaries of the range, whereas arity K is used to determine the number of bins. Each bin is associated with a distinct discrete value. The width of interval is computed by

$$Int_{width} = \frac{max_{val} - min_{val}}{K}.$$
 (6)

Example 2: We continue from Example 1. Consider the example data of numerical type attribute in Table II. In this table, *Age* is a numerical attribute, whereas *Loan approval* is an A^{cls} attribute. Data providers DP₁ and DP₃ own raw data tables Table II(a) and (b), respectively. DP₃ has somewhat different values on the *Age* attribute in contrast to DP₁ on records $\{\text{ID}\#1, 3, 4, 8, 9, 11, 12\}$. They discretize their data on the *Age* attribute, as shown in Table II(c), according to the parameters of equal width binning. A boundary is defined as min_{val} = 10.0 and max_{val} = 70.0, whereas arity K = 5. Though they have differences in their raw data, the produced discrete version is the same for both, since the data values occurred in the specified range. Therefore, the computed information gain 0.34573 is also the same.

Example 3: We continue from Example 1. Consider the raw data tables of two data providers who own common attribute, e.g., *Sex* (which has two values, M or F), as shown in the compressed Table III. The class attribute *Loan approval* shared between the data providers has two values, Y or N, indicating whether or not the loan is approved. Both DP₁ and DP₂ have the same number of records and the same count on their records, i.e., M = 8, and F = 4, but they have different information gain DP₁ = 0.011580 and DP₂ = 0.251629 on the *Sex* attribute. Since the data providers are not consistent in providing the same information on the common RecIDs, this results in a change in the count for class label values. For instance, DP₁ indicates that there is one female whose loan is approved, whereas DP₂ indicates 0 females.

Data Provider DP_1					
ID	Age	Loan approval			
1	39	N			
2	50	N			
3	38	N			
4	53	N			
5	28	N			
6	37	N			
7	49	N			
8	59	N			
9	31	Y			
10	42	Y			
11	37	Y			
12	30	Y			

TABLE II EXAMPLE DATA OF NUMERICAL TYPE ATTRIBUTE

Loan approval

Ν

Ν

N

Ν

N

Ν

Ν

Ν

Y

Y

Ŷ

Y

Data Provider DP₃

ID Age

1 45

 $\mathbf{2}$ 3

> 4 57

5

6 7

8

9

10 42

11

1228

50

35

28

37

49

66

27

38

	Discretization							
ID	Age	Loan approval						
1	[34.0 - 46.0]	Ν						
2	[46.0 - 58.0]	Ν						
3	[34.0 - 46.0]	Ν						
4	[46.0 - 58.0]	Ν						
5	[22.0 - 34.0]	Ν						
6	[34.0 - 46.0]	Ν						
7	[46.0 - 58.0]	Ν						
8	[58.0 - 70.0]	Ν						
9	[22.0 - 34.0]	Y						
10	[34.0 - 46.0]	Y						
11	[34.0 - 46.0]	Y						
12	[22.0 - 34.0]	Y						

Raw data table (c)

Raw data table (a)

Raw data table (b)

TABLE III COMPRESSED DATA TABLE FOR CATEGORICAL TYPE ATTRIBUTE

	Data Provider L	P_1
Sex Loan approval #of l		#of Recs.
М	3Y5N	8
F	1Y3N	4
	Total	12
Raw data table (a)		

4) Computation of Trust Score: Intuitively, the trust score is a metric for assessing the trustworthiness of each data provider. We compute the trust score $TS_{DP_{i}}$ locally for each data provider in an iterative manner on each attribute $\operatorname{Req} A_x$ from the CS. γ is a user-defined weight. A data provider qualifying on the majority gains a positive γ weight in the trust score. On the other hand, a disqualified data provider is penalized with a negative $-\gamma$ weight in the trust score. We aggregate on both positive and negative weights at each iteration to determine the final trust score for each data provider

$$\mathrm{TS}_{\mathrm{DP}_{i}} = \sum_{\mathrm{Req}A_{x} \in \mathrm{Req}A} \gamma \begin{cases} \mathrm{if}(C\mathrm{and} \in \mathrm{Maj}_{C\mathrm{and}}) + \gamma \\ \mathrm{if}(C\mathrm{and} \notin \mathrm{Maj}_{C\mathrm{and}}) - \gamma \end{cases}$$
(7)

B. Security Properties

In this section, we discuss the security properties of our proposed algorithm.

1) Security Against Covert Adversaries: In the context of our problem, a dishonest data provider is a kind of covert adversary who may arbitrarily provide false data on his attribute $A_{\mathcal{T}} \in$ $\mathcal{P}\mathcal{A}_i$. The probability of detecting this cheat by our proposed trust computation algorithm is $1 - \xi$ (refer to the Section V-C1 for details). Each DP_i who has committed to, when registering, the available attributes $\mathcal{PA}_i = \{A_1, \ldots, A_d\}$ is responsible to answer the CS's challenge request, where $\exists \operatorname{Req} A_{x'} \in \mathcal{P} \mathcal{A}_{i}$. When the CS detects a data provider cheating, the provider is penalized with a negative $-\gamma$ weight in the trust score.

2) Mutual Authentication: Before the verification process, each DP_i and the CS mutually authenticate each other by the TLS 1.2 protocol or higher [18], [65]. It is indispensable for the CS to negotiate on the latest stable version of the TLS protocol and stronger cipher suite to prevent against different forms of deception. After successful authentication of each DP_i , they are granted access to the resource queue, where they can register their data attributes.

3) Minimal Access for Outsourcing Verification: The data providers who own customers' private data outsource the verification on their data to the CS. Each DP_i computes locally the information gain function \mathcal{G} on an available attribute $A_{\mathcal{J}} \in \mathcal{P}\mathcal{A}_i$, whereas the CS can have access to only an encrypted $\mathcal{G}'_{A_{\tau}}$ message, i.e., ψ , and its keyed hash, i.e., Υ for the verification. It benefits the data providers to restrict the CS from accessing the customers' private data. Since encrypted individual data records are not exchanged during the verification, the overhead of computation on the CS is also reduced.

4) Authentication and Integrity: HMAC enforces integrity and authenticity. It depends on what underlying hashing function has been used. There are some collision-related vulnerabilities of MD5; however, HMAC-MD5 is not as affected by those vulnerabilities. Regardless, SHA-2 is cryptographically stronger than MD5 and SHA-1. HMAC is constructed by using two nested keys, say k_{in} and k_{out} . These nested keys are not independent; instead, they are derived from a single k_h . Let \mathcal{M} bytes be assumed to be the message blocks for the underlying Merkle–Damgard hash. To derive the keys k_{in} and k_{out} , which are byte strings of length \mathcal{M} , we first construct k_h exactly \mathcal{M} bytes long. If the length of $k_h \leq \mathcal{M}$, we pad it out with zero bytes; otherwise, we replace it with $\mathcal{H}(k_h)$ padded with zero bytes. Then, we compute

$$k_{\text{in}} \leftarrow k_h \oplus \text{ipad}$$

 $k_{\text{out}} \leftarrow k_h \oplus \text{opad}$

The ipad denotes the *inner* pad and the opad denotes the *outer* pad. These pads are 512 bit constants that never change and are embedded in the implementation of the HMAC. The HMAC is

Algorithm 1: IEB_Trust.

KeySetup: CS and DP_i derive *n* symmetric keys by the mutual authentication protocol

Input : Data consumer attributes request

 $\operatorname{Req}A_1, \ldots, \operatorname{Req}A_m$, privacy budget ϵ

Input : Data provider DP_n 's attributes A_1, \ldots, A_d

- Output : Accepted DP_n
- 1: DP_1, \ldots, DP_n own private data tables D_1, \ldots, D_n $\forall i \in I$, where $I = 1, \ldots, n$;
- 2: Each DP_i holds set of attributes $\mathcal{PA}_i = \{A_1, \dots, A_d\}$, over a domain of attributes request Req $A = \{\text{Req}A_1, \dots, \text{Req}A_m\}$;
- 3: $\text{TS}_{\text{DP}_i} \leftarrow 0$; /* Initially, trust score is set to 0 for each data provider */
- 4: $s_t \leftarrow 0$; /* Initially, arrival sequence is set to 0 for all data providers' attributes */
- 5: $\epsilon' \leftarrow \frac{\epsilon}{|\text{Reg}A|};$
- 6: while $\exists \operatorname{Req} A_x \in \operatorname{Req} A$ do
- 7: for $i \in I$ do
- 8: **if** $\exists \operatorname{Req} A_x \in \mathcal{P} \mathcal{A}_i$ then
- 9: register arrival sequence s_t on each attribute;
- 10: end if
- 11: end for
- 12: end while
 - Round 1
- 13: while $\exists \operatorname{Req} A_x \in \operatorname{Req} A$ do
- 14: CS randomly picks $\operatorname{Req} A_x'$ over a range of $\operatorname{Req} A_1, \ldots, \operatorname{Req} A_m$ without replacement;
- 15: CS sends challenge $E(ks_i, \operatorname{Req} A_x')$ to each DP_i where $\exists \operatorname{Req} A_x' \in \mathcal{P} \mathcal{A}_i$;
- 16: Each DP_i computes $\mathcal{G}_{A_{\mathcal{J}}}^{(1)}$ according to (4) and then adds Lap $(1/\epsilon')$, to have $\mathcal{G}_{A_{\mathcal{J}}}^{(1)'}$;
- 17: Each DP_i encrypts the message $\psi^{(1)} \leftarrow E(ks_i, \mathcal{G}_{A_{\mathcal{J}}}^{(1)'})$ and then computes tag $\Upsilon^{(1)} \leftarrow \mathcal{S}(k_h, \psi^{(1)});$
- 18: CS receives $\psi^{(1)} \| \Upsilon^{(1)} \| DP_i$ on his challenge from the corresponding data providers;
- 19: CS computes comparison to determine $\operatorname{Maj}_{Cand}^{R(1)} \leftarrow$ findMajCand($\psi^{(1)} \| \Upsilon^{(1)}$, size);
- 20: end while Round 2
- 21: while $\exists \operatorname{Req} A_x \in \operatorname{Req} A$ do
- 22: **for** $\ell = 1$ to α **do**
- 23: CS generates \mathcal{K} random IDs for Req A_x' (pick in Round 1) from $|D_i|$ records, where $5 \le \mathcal{K} \le 10$;
- 24: CS generates \mathcal{P} pairs of values for Req A_x' and A^{cls} attributes, where $5 \leq \mathcal{P} \leq 10$;
- 25: CS sends challenge $E(ks_i, \mathcal{K}, \operatorname{Req} A_{x'}) || E(ks_i, v_{x'}, v_{\operatorname{cls}})$ to each DP_i where $\exists \operatorname{Req} A_{x'} \in \mathcal{P} \mathcal{A}_i$;
- 26: Each DP_i computes $\mathcal{G}_{A_{\mathcal{J}}}^{(2)}$ on the concatenated \mathcal{K} specified records and \mathcal{P} pairs of values and then adds Lap $(1/\epsilon')$, to have $\mathcal{G}_{A_{\mathcal{J}}}^{(2)'}$;

- 27: Each DP_i encrypts the message $\psi^{(2)} \leftarrow E(ks_i, \mathcal{G}_{A_{\mathcal{J}}}^{(2)'})$ and then computes tag $\Upsilon^{(2)} \leftarrow \mathcal{S}(k_h, \psi^{(2)});$
- 28: CS receives $\psi^{(2)} \| \Upsilon^{(2)} \| DP_i$ on his challenge from the corresponding data providers;
- 29: CS computes comparison to determine $\operatorname{Maj}_{Cand_{\ell}}^{R(2)}$ $\leftarrow \operatorname{find}\operatorname{Maj}Cand(\psi^{(2)} || \Upsilon^{(2)}, \operatorname{size});$
- 30: end for
- 31: CS computes $\bigcap_{\ell=1}^{\alpha} \operatorname{Maj}_{Cand_{\ell}}^{R(2)}$;
- 32: end while
- 33: CS computes $\operatorname{Maj}_{Cand}^{R(1)} \cap \operatorname{Maj}_{Cand}^{R(2)}$ to determine $\operatorname{Maj}_{Cand}$;
- 34: **for all** C and \in Maj_{Cand} **do**
- 35: set *C* and as $\text{Qual}_{\text{DP}_i}$;
- 36: $TS_{DP_i} = TS_{DP_i} + \gamma;$
- 37: end for
- 38: for all C and \notin Maj_{Cand} do
- 39: set C and as UnQual_{DP_i};
- 40: $TS_{DP_i} = TS_{DP_i} \gamma;$
- 41: **end for**
- 42: if size == $1 \wedge TS_{DP_i} \ge 0$ then
- 43: set DP_i as $Qual_{DP_i}$;
- 44: **end if**
- 45: Pick one *C* and by comparison on the arrival sequences of the $Qual_{DP_n}$ on each attribute;
- 46: **return** Data providers whose final aggregated trust score ≥ 0

assumed to be a secure PRF [14]. It provides better protection against length extension attacks. It is built as follows:

$$\mathcal{S}(k_h, \psi) = \mathcal{H}(k_h \oplus \text{opad}, \mathcal{H}(k_h \oplus \text{ipad} \| \psi)).$$

One of the properties of a cryptographic hash function is that if there is a minor change in an input message, it changes the message digest so extensively that the new message digest appears uncorrelated with the old computed message digest. In our case, we do not apply cryptographic hash functions directly on the input data for data integrity because we allow parties to have minor inaccuracies on numerical attributes for a specified threshold.

C. Analysis

In this section, we analyze the correctness and security of Algorithm 1.

Proposition V.1 (Correctness): Assuming multiple data providers are dishonest, Algorithm 1 correctly computes the trust scores among them, as stated in Problem 1 in Section II-B, to evaluate the trustworthiness of each data provider.

Proof: Algorithm 1 selects an attribute uniformly at random without replacement from a list $\operatorname{Req} A = {\operatorname{Req} A_1, \ldots, \operatorname{Req} A_m}$ of *m* requested attributes. Each DP_i computes $\mathcal{G}_{A_{\mathcal{J}}}$ according to (4) for its matching attribute in the presence of a shared class attribute A^{cls} . For a continuous-valued attribute, each provider follows equal-width method for discretization into intervals

of equal size. Consider $A_{\mathcal{T}}$ is discrete-valued, owned by two providers, where $\Omega(A_{\mathcal{J}}) = \{v_1, v_2\}$ is in its domain of data values. Assume that there is a single record between two providers, where they have different values. Algorithm 1 computes $\mathcal{G}_{A_{\mathcal{J}}}^{(1)'}$ in the first round for both the data providers and returns different scores. This suggests that they are not the same.

Now, we consider an extended case, where two data providers (say DP_1 and DP_2) would have different sets of records, but the computation of $\mathcal{G}_{A_{\mathcal{J}}}^{(1)'}$ in the first round on the full dataset for both data providers returns the same score, so we have $\operatorname{Maj}_{Cand}^{R(1)} = \{\operatorname{DP}_1, \operatorname{DP}_2\}.$ Algorithm 1 verifies further by running the process α times in the second round. During each iteration, data providers have to select records over \mathcal{K} random IDs for $A_{\mathcal{I}}$, and they also have to add \mathcal{P} pairs of values $v_{x'}$ and v_{cls} for $A_{\mathcal{J}}$ and a class attribute A^{cls} , respectively, from the CS before computa-tion of $\mathcal{G}_{A_{\mathcal{I}}}^{(2)'}$. Algorithm 1 computes $\operatorname{Maj}_{Cand}^{R(1)} \cap (\bigcap_{\ell=1}^{\alpha} \operatorname{Maj}_{Cand_{\ell}}^{R(2)})$ to determine $\operatorname{Maj}_{Cand}$. This determines whether or not the data providers are holding the same data values over the common attribute $A_{\mathcal{J}}$. Data providers are required to match in both the rounds to prove that they have the same score. Since data providers are holding a different set of records, it is not possible for them to match because of the randomness introduced in the second round.

Proposition V.2 (Security): Algorithm 1 is secure against covert adversaries, as described in Section V-B1, by the probabilistic bound of $1 - \xi$.

Proof: The security of Algorithm 1 depends on the key derivation in the mutual authentication protocol and the communication of the CS and data providers DP_n in the verification process.

- 1) A random challenge $E(ks_i, \text{Req}A_x')$ is secure because of symmetric keys derivation in [18] and [65].
- 2) On a given challenge request, if $\exists \operatorname{Req} A_x' \in \mathcal{P} A_i$, each data provider first computes the information gain function on its matching attribute $\mathcal{G}_{A_{\mathcal{J}}} \in \mathcal{PA}_i$ and then perturbs the output by adding noise. This returns a noisy score $\mathcal{G}'_{A_{\tau}}$, for which data providers should agree on the scale for digits after the decimal point. It is secured for privacy protection because each DP_i only exchanges an encrypted $\mathcal{G}'_{A_{\mathcal{T}}}$ message, i.e., ψ , and its keyed hash, i.e., Υ , with the CS in both rounds of the protocol, instead of exchanging encrypted individual data records on their attributes $A_{\mathcal{J}}$.
- 3) Keyed HMAC $S(k_h, \psi)$ is a secure PRF according to [14]. It is computationally infeasible for an adversary to find distinct inputs ψ_1, ψ_2 such that $S(k_h, \psi_1) = S(k_h, \psi_2)$.
- 4) Dishonest data providers cannot modify the outputs, i.e., $\psi \| \Upsilon$, of the honest providers in any round of the protocol. They may compute $\mathcal{G}^*_{A_{\mathcal{I}}}$ on their false data and can send their $\psi^* \| \Upsilon^*$ to the CS. The CS computes a comparison and detects cheating from a dishonest data provider with the probability of $1 - \xi$.

1) Adversary's Inferences: In the following, we estimate the probability of an adversary, i.e., a dishonest data provider, to correctly guess $\mathcal{G}^*_{A_{\tau}}$ on a random challenge attribute $\operatorname{Req} A'_x$. An adversary knows $|D^{\tau}|$, the number of records in D^{τ} , and $|A_{i,D^{\tau}}^{cls}|$, the number of records of class A_i^{cls} in D^{τ} , and computes the entropy of D^{τ} by (2). Next, the adversary may try to compute entropy on $A_{\mathcal{T}}$ by the following equation because he knows $|\Omega(A_{\mathcal{T}})|$, the domain size of $A_{\mathcal{T}}$, and $|D^{\tau}|$, the number of records in D^{τ}

$$E_{A_{\mathcal{J}}}^{*}(D^{\tau}) = \sum_{j'=1}^{\mathcal{V}'} \frac{|D_{j'}^{\tau}|}{|D^{\tau}|} \times -\sum_{i=1}^{\mathcal{C}} \frac{|A_{i,D_{j'}}^{\text{cls}}|}{|D_{j'}^{\tau}|} \times \log_2 \frac{|A_{i,D_{j'}}^{\text{cls}}|}{|D_{j'}^{\tau}|}.$$
(8)

There are $|\Omega(A_{\mathcal{J}})|^{|D'|}$ possible arrangements, in which an adversary may try to compute $E^*_{A_{\tau}}(D^{\tau})$. Finally, he computes $\mathcal{G}^*_{A_{\sigma}}$ having all distinct values by the following equation:

$$\mathcal{G}_{A_{\tau}}^{*} = E(D^{\tau}) - E_{A_{\tau}}^{*}(D^{\tau}).$$
(9)

This results in ϑ distinct values of $\mathcal{G}^*_{A_{\tau}}$, with the lower bound of $\vartheta \approx |D^{\tau}|$. The probability of correctly guessing $\mathcal{G}^*_{A_{\tau}}$ for an adversary in our verification process is

$$\xi = \frac{1}{\vartheta} \times \left(\frac{1}{\vartheta}\right)^{\alpha}.$$
 (10)

2) Detecting Cheat Against Varying Dishonest Providers: Let n denote the number of participating data providers, and let b denote an upper bound on the number of dishonest data providers who may arbitrarily provide incorrect data in responding to the CS's challenge.

- 1) When b < n/2, the verification process guarantees fairness, and no honest data providers are negatively affected by their trust levels.
- 2) When $b \le n-2$, the verification process guarantees fairness under the arbitrary behavior of dishonest data providers, where the chance of detecting them is $1 - \xi$. It is a type of covert adversarial behavior when the dishonest data providers arbitrarily provide false data on their data inputs, i.e., they neither would be able to appear in the majority nor would be able to undermine the reputations of the honest data providers.
- 3) When b > n/2, the verification process does not guarantee fairness on the flip side, i.e., when the behavior of dishonest data providers is not arbitrary. This would be the case when the dishonest data providers not only appear in the majority, but also organize in a way to undermine the reputation of the honest data providers. We assume that if a secure set intersection is carried out by using a trusted mediator (e.g., by computing the function on the data providers input) between data providers, then the dishonest providers would not be able to determine the total number of participating data providers in advance. This would restrict them from developing the organized group; still, there is no remedy if they would try by guessing at random.

D. Evaluation of Learner Models

We provide an example of a sample data to evaluate the quality of linear regression, kNN, and random forest learner models.

Example 4: We retrieve the top 1000 records from a real-life Adult¹ dataset on attributes age, education-num, race, sex, and

¹Available at: http://archive.ics.uci.edu/ml/datasets/Adult

income. The attributes *age and education-num* are of continuous types, whereas *race, sex, and income* are of categorical types. We develop learner models in *RapidMiner*² to compare the predictive accuracy of linear regression, kNN, and random forest methods.

For the linear regression model, we set *education-num* as a label, which is considered as a dependent attribute (or variable), and the remaining are considered as independent attributes. We convert nonnumeric type attributes to the numeric type. After running tenfold cross validation, the *root-mean-square error* is found to be 2.438 ± 0.165 , which indicates the standard deviation of the residuals. Furthermore, R^2 is found to be 0.127 ± 0.055 , which indicates the goodness of fit of this regression model. Its value is close to 0, indicating a weak linear correlation.

For the kNN model, we set all attributes as nominal and *education-num* as a label. After running tenfold cross validation when k = 20, the *accuracy* is found to be 33.90% ± 5.59%, which indicates the percentage of correct predictions.

For the random forest model, we set the *education-num* attribute as nominal and specify the role as a label. The key parameter "number of trees" is specified as 10, and the "gain ratio" is chosen as a criterion for splitting attributes. After running tenfold cross validation, the *accuracy* is found to be $32.90\% \pm 0.30\%$, which indicates the percentage of correct predictions.

There are no significant performance differences found on running these learner models on the sample dataset. Data providers would use any one or multiple learning methods for missing data imputation.

E. Price Setting Using Auction Mechanism

An auction mechanism can be defined in many different ways depending upon the design requirements. The two variants of second price sealed-bid auctions [23] have been widely used, namely, VCG and generalized second price (GSP) mechanisms for multiple items.

The reason for employing the VCG mechanism for determining the pricing on data providers' attributes is that truthful bidding is a dominant strategy, and there is no incentive to lie or deviate from reporting true valuations for a data provider. It maximizes the total valuation obtained by data providers. One nice property of the VCG mechanism is that it provides a unique outcome, which is socially optimal, whereas, in the GSP, there would be multiple outcomes in terms of Nash equilibrium. One Nash equilibrium would maximize social welfare but not all of them.

We intend to design an auction mechanism for multiple items. It is assumed that the data providers intend to set up a matching market using a second price sealed-bid auction for valuation of their attributes. We formally define the procedure for setting the price as follows.

1) Data Providers: Let DP_1, \ldots, DP_n (where $i = 1, \ldots, n$) be the set of data providers, who set up a matching market for valuations of their attributes.

2) Positions: Let P_1, \ldots, P_n (where $j = 1, \ldots, n$) be the set of positions for which data providers compete. The higher the position P_j , the more will be its demand rate. The positions should be equal to the number of data providers. If there are more data providers than positions, we simply add fictitious positions of demand rate 0. Similarly, if there are more positions than data providers, we add fictitious data providers of revenue per demand 0.

3) Revenue Per Demand: Revenue per demand is the expected amount of money that a data provider DP_i receives, denoted by Rev_i , for every demand on its attribute. The monetary values of Rev_i are sorted in descending order.

4) Demand Rate: The demand rate is defined as the number of demands requested by a consumer over a period of time, denoted by Q_j . The demand rate varies as per the position P_j . Q_j enumerates in descending order.

5) Data Providers' Valuations: Data providers' valuations are defined as the data provider DP_i 's valuation of the position P_j . It is the product of the revenue per demand Rev_i and the demand rate Q_j , denoted by $Val_{i,j}$. It is computed as follows:

$$\operatorname{Val}_{i,i} = \operatorname{Rev}_i \times Q_i. \tag{11}$$

6) VCG Price: VCG price is defined as the harm or externality caused by data provider DP_i to other data providers in terms of reduction of their valuations due to his presence. It is called VCG price, denoted by $ExPrc_{i,j}$, which is paid by data provider DP_i for position P_j . Formally, it is defined by

$$\operatorname{ExPrc}_{i,j} = \bigvee_{\operatorname{DP}_n - \operatorname{DP}_i}^{P_n} - \bigvee_{\operatorname{DP}_n - \operatorname{DP}_i}^{P_n - P_j}$$
(12)

where

- 1) $DP_n DP_i$ is the set of data providers excluding data provider DP_i ;
- 2) $P_n P_j$ is the set of positions excluding position P_j ;
- 3) $\bigvee_{\text{DP}_n-\text{DP}_i}^{P_n}$ is the sum of data provider values of an optimal matching between sets $\text{DP}_n \text{DP}_i$ and P_n ; and
- 4) $\bigvee_{\text{DP}_n-\text{DP}_i}^{P_n-P_j}$ is the sum of data provider values of an optimal matching between sets $\text{DP}_n \text{DP}_i$ and $P_n P_j$.

7) Data Providers' Valuations After Payoff: Data providers' valuations after payoff is defined as the data provider DP_i 's valuation on position P_j after paying off harm to other data providers. It is calculated using the following equation:

$$Val_{DP_i} = \max Val_{i,j} - ExPrc_{i,j}.$$
 (13)

8) Valuation of an Attribute: Valuation of an attribute can be assessed once a data provider DP_i acquires a certain position P_j . The value of each data provider's attribute per single demand is calculated using the following equation:

$$ValAttr_{DP_i} = \frac{Val_{DP_i}}{Q_j}.$$
 (14)

9) Attribute Count: The attribute count $CntAttr_{DP_i}$ of a data provider DP_i represents the number of attributes in a single record. Each DP_i owns a mutually exclusive set of attributes.

10) Price Per Record: The price per record $PrcRec_{DP_i}$ of a data provider DP_i represents the unit price of a record. Naturally,

²Available at: https://rapidminer.com/products/studio/

it is the product of the value per attribute $ValAttr_{DP_i}$ and the attribute count CntAttr_{DP_i} in a single record. That is,

$$\operatorname{PrcRec}_{\operatorname{DP}_{i}} = \operatorname{ValAttr}_{\operatorname{DP}_{i}} \times \operatorname{CntAttr}_{\operatorname{DP}_{i}}.$$
 (15)

11) Size of the Dataset: The dataset of each data provider DP_i consists of a collection of records, denoted by $|D_i|$. The size of a dataset grows as the number of records in the dataset increases.

12) Price of the Raw Dataset: The price of a raw dataset $PrcRawDS_{DP_i}$ represents the data provider DP_i 's selling price of a raw dataset in the e-market. The overall pricing of a raw dataset increases as the number of records or the unit *price per record* increases. It is computed as follows:

$$\operatorname{PrcRawDS}_{\operatorname{DP}_{i}} = |D_{i}| \times \operatorname{PrcRec}_{\operatorname{DP}_{i}}.$$
 (16)

13) Total Price of the Raw Dataset: The total price of the raw dataset $\text{TPrc}_{\text{RawDS}}$ is the sum of the pricing of all the contributing data providers' raw datasets. It is computed as follows:

$$TPrc_{RawDS} = \sum_{i=1}^{n} PrcRawDS_{DP_i}.$$
 (17)

14) Total Monetary Value of the Raw Dataset: First, data providers compute baseline accuracy (BA) for classification analysis using the secure multiple party classifier [21] by maintaining the confidentiality of their raw data. Then, they use the information utility of classifying raw data to derive the monetary value of the raw dataset, denoted by TMValue_{RawDS}. It is calculated using the following equation:

$$TMValue_{RawDS} = TPrc_{RawDS} \times BA.$$
(18)

F. Anonymization Method

In this section, we provide an extension of the two-party *Differentially* private anonymization in Algorithm 2, which is based on *Generalization* [53] to differentially integrate multiple private data tables. This algorithm guarantees ϵ -differential privacy and security definition under the semihonest adversary model (readers may refer to the detailed analysis in [53, sec. 6.3]). The two major extensions over the TDS algorithm [31] include: 1) *DistDiffGen* selects the *Best* specialization based on the exponential mechanism, and 2) *DistDiffGen* perturbs the generalized contingency table by adding the Laplacian noise to the count of each equivalence group.

Generally, there is no incentive for any data provider who executes the algorithm, as the purpose is merely to synchronize the anonymization process. We assume that a trusted data provider, who attains the highest trust score after running Algorithm 1, starts the anonymization process. The accepted data providers, as a result of trust computation by Algorithm 1, attain a mutually exclusive set of attributes, i.e., $\mathcal{PA}_i \cap \mathcal{PA}_j = \emptyset$ for any $1 \le i, j \le n$ over the same set of records for integrating data.

Initially, all values in the set of attributes $\mathcal{PA}_i = \{A_1, \ldots, A_d\}$ of each data provider are generalized to the topmost value in their taxonomy trees (Line 1), as illustrated in Fig. 1, and Mark_{κ} contains the topmost value for each attribute $A_{\mathcal{J}} \in \mathcal{PA}_i$ (Line 2). Each data provider keeps a copy of the \cup Mark_{κ} and a generalized data table D_g . The attribute $A_{\mathcal{J}}$ can be either categorical or numerical, but the class attribute is required to be categorical. The split value of a categorical attribute v_c is a generalized value drawn from a predefined taxonomy tree of the attribute, whereas the split value of a numerical attribute v_{num} is determined by using the exponential mechanism (Line 4). It partitions the domain range of a numerical attribute into successive intervals $\mathcal{I}_1, \ldots, \mathcal{I}_k$. Line 4 preserves $\epsilon' |A_{num}|$ -differential privacy since the cost of each exponential mechanism is ϵ' . In Line 5, a score IGScore is computed for all candidates $v \in \bigcup Mark_{\kappa}$. At each iteration, the algorithm uses the secure distributed exponential mechanism (DistExp) as presented in [53] (readers may refer to the details of the DistExp algorithm) to select a winner candidate $w \in$ \cup Mark_{κ} for specialization (Line 7). Different utility functions (e.g., information gain) can be used to calculate the score. If the winner candidate w is local to DP_i , DP_i specializes w on D_q by splitting its records into child partitions, updates its local copy of \cup Mark_{κ}, and instructs all the other participating data providers to specialize and update their local copy of \cup Mark_{κ} (Line 8–11). The information gain, denoted by \mathcal{G}_{DP_i} , accumulates IGScore(x) on the winner's attribute specializations (Line 12). DP_i further calculates the scores of the new candidates as a result of the specialization (Line 14). If the winner w is not one of DP_i 's candidates, DP_i waits for instructions from the other winner data provider DP_i, where $i \neq j$, to specialize w and to update its local copy of \cup Mark_{κ} (Lines 16 and 17). This process iterates until the specified number of the specializations h is reached. The algorithm perturbs the output by adding the noisy count at each leaf node (Line 21) using the Laplace mechanism. The contribution of each data provider is computed according to (22). Finally, the monetary share of each data provider is derived according to the (23).

G. Quantifying the Monetary Value

The rationality of quantifying the monetary value is that data providers are the business stakeholders, who collaborate in the data integration process to maximize their profits. The profit generated by their collaboration is distributed based on each provider's contribution to information utility and its trustworthiness.

1) Cost of Anonymization in Integrated Data: First, the data providers compute classification accuracy (CA) on the anonymized integrated data. Then, they quantify the cost of anonymization in integrated data, denoted by Cost_{IntDS}, on the difference between the BA and the CA. It is computed as follows:

$$Cost_{IntDS} = TPrc_{RawDS} \times (BA - CA).$$
 (19)

2) Expected Value in Integrated Data: An expected monetary value in integrated data is what the data providers earn from the information utility of classification analysis when trading an anonymized version of integrated data. The information utility varies with the valuations of data providers' attributes and joint privacy requirements, such as privacy budget ϵ and specialization level *h*, for a ϵ -differential privacy model in a distributed setup,

Algorithm 2: Monetary Shares for Data Providers Using DistDiffGen.

Input : Data providers' attributes valuations ValAttr_{DP_n} **Input** : Private data tables D_1, \ldots, D_n , privacy budget ϵ , and number of specializations h

Output : Monetary shares MShare_{DP_n}

- 1: Initialize D_q with one record containing topmost generalized values in each data provider's taxonomy tree;
- 2: Initialize Mark_{κ} to include the topmost value;
- 3:
- $\epsilon' \leftarrow \frac{\epsilon}{2(|A_{\text{num}}|+2h)};$ Determine the split value for each $v_{\text{num}} \in \bigcup \text{Mark}_{\kappa}$ 4: with probability $\propto \exp(\frac{\epsilon'}{2\Delta u}u(D, v_{\text{num}}));$
- Compute the IGScore for $\forall v \in \cup Mark_{\kappa}$; 5:
- for iter = 1 to h do 6:
- 7: Determine the winner candidate w by using the DistExp Algorithm [53];
- 8: if w is local then
- Specialize w on D_g ; 9:
- 10: Replace w with child(w) in the local copy of \cup Mark_{κ};
- 11: Instruct all the other participating data providers to specialize and update \cup Mark_{κ};
- $\mathcal{G}_{\mathrm{DP}_i} = \mathcal{G}_{\mathrm{DP}_i} + \mathrm{IGScore}(x);$ 12:
- 13: Determine the split value for each new $v_{\text{num}} \in \bigcup \text{Mark}_{\kappa}$ with probability \propto $\exp(\frac{\epsilon}{2\Delta u}\mathbf{u}(D, v_{\text{num}}));$
- 14: Compute the IGScore for each new $v \in \bigcup Mark_{\kappa}$;
- 15: else
- 16: Wait for the instruction from the winner data provider:
- 17: Specialize w and update \cup Mark_{κ} using the instruction;

18:
$$\mathcal{G}_{\mathrm{DP}_i} = \mathcal{G}_{\mathrm{DP}_i} + \mathrm{IGScore}(x);$$

- 19: end if
- 20: end for
- 21: Compute count $(CT + Lap(2/\epsilon))$ for each leaf node;
- 22: Compute the contribution of each data provider according to (22);
- 23: Compute monetary share of each data provider according to (23);
- 24: return MShare_{DP}

between the data providers. It is calculated on the difference between the total monetary value of the raw dataset TMValue_{RawDS} and the cost of anonymization in integrated data Cost_{IntDS}. It is computed as follows:

$$EValue_{IntDS} = TMValue_{RawDS} - Cost_{IntDS}.$$
 (20)

3) Expected Value of an Individual Data Provider: The expected monetary value of an individual data provider, denoted by EValueIndv_{DP_i}, is determined by the ratio of the number of attributes $CntAttr_{DP_i}$ a data provider owns with the total count of attributes. It is computed as follows:

$$EValueIndv_{DP_i} = EValue_{IntDS} \times \frac{CntAttr_{DP_i}}{\sum_{i=1}^{n} CntAttr_{DP_i}}.$$
 (21)

4) Derivation of Monetary Share: The derivation of a monetary share depends upon the contribution of each data provider and its trustworthiness. Intuitively, a data provider, whose provided data on his attributes result in more information gain, and whose trust level is higher than the other competitors, can get a significantly larger share of the monetary value. The contribution of each data provider DP_i is derived from the expected monetary value EValueIndv_{DP_i} by fairly computing first the accumulative information gain \mathcal{G}_{DP_i} of each data provider DP_i on the anonymized integrated dataset. The information gain IGScore(x) of the winner candidate w data provider accumulates under the relevant winner w data provider at each iteration (refer to the Section V-F for details) for the specified specialization level h. The contribution of each data provider Contrib_{DP_i} is calculated using the following equation:

$$\text{Contrib}_{\mathsf{DP}_i} = \frac{\widetilde{\mathcal{G}}_{\mathsf{DP}_i}}{\sum_{i=1}^{n} \widetilde{\mathcal{G}}_{\mathsf{DP}_i}} \times \text{EValueIndv}_{\mathsf{DP}_i}.$$
 (22)

Finally, the monetary share of each data provider MShare_{DP}, is derived according to (7), i.e., the aggregated trust score of each data provider, and (22), i.e., the contribution of each data provider. Therefore, MShare_{DP}, becomes

$$\text{MShare}_{\text{DP}_i} = \text{Contrib}_{\text{DP}_i} \left(1 + \frac{\text{TS}_{\text{DP}_i}}{\sum_{i=1}^n \text{TS}_{\text{DP}_i}} \right).$$
(23)

VI. COMPARATIVE ANALYSIS AND EMPIRICAL STUDY

In this section, we first provide a comparison of our approach, followed by an empirical study.

A. Comparative Analysis

We compare our proposed IEB_Trust, an entropy-based trust computation algorithm with the closely related provenancebased trust method [16]. The provenance-based method computes the trust scores for data and data providers using similarity functions, but do not consider privacy protection when evaluating trustworthiness. The fundamental idea of our approach is different. Our method enables secure trustworthiness assessment and preserves the privacy of the customers' data when evaluating the trustworthiness of the participating data providers. For this reason, we are limiting to the runtime comparison in Fig. 3(a). We evaluate the performance of our proposed method on a real-life Adult³ dataset. It contains 45 222 records with eight categorical attributes, six numerical attributes, and a binary class attribute Income with two levels, $\leq 50K$ or > 50K. The distribution of attributes other than class attribute among ten data providers is shown in Fig. 3(b). We generate 10% of data conflicts over randomly chosen attributes. We vary the size of the datasets $|D_i|$ from 10 to 50K to study the runtime cost.

³Available at: http://archive.ics.uci.edu/ml/datasets/Adult



Fig. 3. Our method improves the runtime efficiency compared to the provenance-based trust method. (a) Runtime comparison. (b) Distribution of attributes.

All experiments are conducted on an Intel Core i7 3.4 GHz PC with 8-GB memory.

The running time includes time elapsed in both the initialization phase and the iteration phase. We observe that the initialization phase of the provenance-based method takes more time to compute data similarity and data conflict. It has worst-case complexity of $O(n^2)$, while the complexity of our proposed method at the initialization phase is $O(\text{CntAttr}_D P_i \cdot |D_i| \log |D_i|)$. Since each data provider computes $\mathcal{G}_{A_{\mathcal{J}}}$ in a distributed setup, the complexity remains the same in our method. The iteration phase to compute trust is much faster in both the methods. It takes less than 1 s to complete the trust computation. Fig. 3(a) shows that our method is more efficient in running time over the provenance-based method. Our method is scalable when we need to grow either the number of attributes, the number of data providers, or both on a dataset.

B. Empirical Study

We first analyze the trustworthiness of each data provider and assess the truthfulness of the provided data by a trust score metric. Second, we analyze the impact of ϵ -differential privacy requirements along with the aggregated trust score on each data provider's monetary value. We evaluate our proposed method, *IEB_Trust*, with the assumption of having four data providers who intend to verify the correctness of their data before participation in the data mashup. This assumption is reasonable because we have a limited number of attributes in the dataset to be shared among data providers.

1) Trust Measurement: Our proposed method evaluates the trust of participating data providers based on the following conditions.

- 1) A data provider is found as honest and gains a positive score.
- A data provider is found as dishonest and is penalized with a negative score.

- 3) A single data provider of an attribute that no others own is accepted based on the existing trust score $TS_{DP_i} \ge 0$ without an increase in the trust score.
- 4) A data provider who does not register for an attribute has no effects on the trust score.

To demonstrate the effectiveness of our approach, we conduct two cases of experiments that are independent of each other. This means that for each case, data providers hold different sets of overlapping attributes with their arrival sequences. In each case, we assume $\gamma = 0.5$, but it does not need to be fixed to a specific weight.

Consider the first case with the participating data providers' attributes and their arrival sequences. $DP_1 \mapsto A_1:s_{t_1}, A_7:s_{t_1}, A_8:s_{t_1}, A_9:s_{t_1}, A_{10}:s_{t_2}, A_{11}:s_{t_1}; DP_2 \mapsto A_2:s_{t_2}, A_3:s_{t_1}, A_4:s_{t_1}, A_5:s_{t_2}, A_7:s_{t_2}, A_8:s_{t_3}, A_{13}:s_{t_1}; DP_3 \mapsto A_1:s_{t_2}, A_4:s_{t_2}, A_5:s_{t_1}, A_6:s_{t_1}, A_8:s_{t_2}, A_{11}:s_{t_2}, A_{13}:s_{t_2}; and DP_4 \mapsto A_1:s_{t_3}, A_2:s_{t_1}, A_5:s_{t_3}, A_9:s_{t_2}, A_{10}:s_{t_1}, A_{11}:s_{t_3}, A_{12}:s_{t_1}.$ Fig. 4(a) depicts the trust scores analysis for Case 1 based on the demand of a data consumer on attributes A_1, \ldots, A_{13} .

It is observed that the DP₂ trust score never drops during the verification process in contrast to the other competing data providers. The flat lines from A_2 to A_6 at trust score level 0.5, and A_9 to A_{12} at trust score level 2.5, indicate that those attributes are not submitted by DP_1 and DP_2 , respectively. This is not always the case; for instance, there are flat lines from A_2 to A_3 at trust score level 0.5, A_5 to A_6 at trust score level 0.5, and A_{11} to A_{12} at trust score level 2.0, indicating that DP₂, DP₃, and DP₄ are the single data providers on those attributes. DP2, DP3, and DP₄ are accepted because they are maintaining an aggregated trust score ≥ 0 at that point of the verification. However, their trust scores do not increase because they own an attribute that no others own. It is assumed that DP_1 has 5% of missing data on A_8 and A_{11} , DP₃ has 5% of missing data on A_5 , and DP₄ has 1% of missing data on A_1 . They impute missing data by using the kNN imputation method in order to claim it as original data. Our trust verification approach restricts this dishonest behavior of data providers; for instance, DP₁ at A_8 and A_{11} , DP₃ at A_5 , and DP₄



Fig. 4. Trust scores analysis. (a) Case 1. (b) Case 2.



Fig. 5. Aggregated trust scores. (a) Case 1. (b) Case 2.

at A_1 , by penalizing them with negative weight in their trust scores. Fig. 5(a) depicts the aggregated trust scores for Case 1. DP₂ attains the maximum trust score 3.0 in competing with the other data providers, whereas DP₁ ends up with the minimum trust score 1.0. There is a tie on aggregated trust scores between DP₃ and DP₄.

Consider the second case with the participating data providers' attributes and their arrival sequences. $DP_1 \mapsto A_1:s_{t_1}$, $A_6:s_{t_3}$, $A_7:s_{t_1}$, $A_8:s_{t_2}$, $A_9:s_{t_3}$, $A_{10}:s_{t_2}$, $A_{12}:s_{t_2}$; $DP_2 \mapsto$ $A_2:s_{t_2}$, $A_5:s_{t_2}$, $A_6:s_{t_4}$, $A_7:s_{t_2}$, $A_8:s_{t_1}$, $A_9:s_{t_2}$, $A_{11}:s_{t_1}$; $DP_3 \mapsto A_3:s_{t_1}$, $A_5:s_{t_1}$, $A_6:s_{t_1}$, $A_8:s_{t_3}$, $A_9:s_{t_1}$, $A_{12}:s_{t_1}$, $A_{13}:s_{t_2}$; and $DP_4 \mapsto A_2:s_{t_1}$, $A_4:s_{t_1}$, $A_6:s_{t_2}$, $A_9:s_{t_4}$, $A_{10}:s_{t_1}$, $A_{11}:s_{t_2}$, $A_{13}:s_{t_1}$. Fig. 4(b) depicts the trust scores analysis for Case 2 based on the demand of a data consumer on attributes A_1, \ldots, A_{13} .

It is observed that DP₁, DP₂, and DP₄ maintain their trust scores quite well except for a fall of 0.5 in their trust scores at A_9 , A_5 , and A_{13} , respectively. The flat lines from A_1 to A_5 at trust score level 0.0, and A_3 to A_5 at trust score level 0.5, indicate that those attributes are not submitted by DP₁ and DP₄, except at A_1 and A_4 , respectively. Since DP₁ and DP₄ are the single data providers on A_1 and A_4 , their trust scores do not increase. However, they are accepted because they maintain an aggregated trust score ≥ 0 . We observe that DP₃ is inconsistent in maintaining its trust level throughout the verification process. It is worthwhile to note that our trust verification process restricts the arbitrary behavior of dishonest DP₁ and DP₃ to undermine the trust levels of DP₂ and DP₄. Fig. 5(b) depicts the aggregated trust scores for Case 2. DP₂ attains the maximum trust score 2.5 in competing with the other data providers, whereas DP₃ ends up with a negative trust score of -1.0. This results in the rejection of DP₃ from the final selection in the data mashup.

2) Impact of Privacy Protection and Trust Score on DP's Monetary Value: In this section, we analyze the impact of ϵ -differential privacy requirements along with the aggregated trust score on each data provider's monetary value. Recall from Section V-E that both revenue per demand Rev_i and demand rate Q_j are enumerated in descending order. Suppose Rev_i = {0.6, 0.5, 0.4, 0.3} and $Q_j = {0, 0.7, 0.6}$ for data providers DP₁, DP₂, DP₃, and DP₄, respectively. The inputs for Rev_i and Q_j do not need to be fixed to a particular value, it is just assumed here for simplicity.

Case 1: Table IV(a) shows the selection of attributes from each accepted data provider. BA on the integrated data of accepted



Fig. 6. Impact of ϵ -differential privacy requirements and trust scores on DP₁, DP₂, DP₃, and DP₄ monetary value (Case 1). (a) $\epsilon = 0.2$. (b) $\epsilon = 0.4$. (c) $\epsilon = 0.6$. (d) $\epsilon = 0.8$.

TABLE IV Selection of Attributes From Data Providers

DP_1	DP_2	DP_3	DP_4		DP_1	DP_2	DP_4
A_1	A_5	A_8	A_2		A_1	A_8	A_2
A_9	A_4	A_6	A_{12}		A_7	A_9	A_4
A7	A ₁₃	A ₁₁	A_{10}		A ₁₂	A ₁₁	A ₁₀
	A_3						A_6
(a) Case 1				-	(b) Case 2		

data providers is 85.3% using the secure multiple-party classifier [21] without disclosing their raw data. We vertically partition the *Adult* dataset into four partitions VP₁, VP₂, VP₃, and VP₄ for data providers DP₁, DP₂, DP₃, and DP₄, respectively. Furthermore, we split the dataset into 30 162, and 15 060 records for the training and testing sets, respectively. The valuation of each data provider's attribute is \$0.47, \$0.41, \$0.36, and \$0.30, representing ValAttr_{DP1}, ValAttr_{DP2}, ValAttr_{DP3}, and ValAttr_{DP4} by (14). The attribute count of each data provider is CntAttr_{DP1} = 3, CntAttr_{DP2} = 4, CntAttr_{DP3} = 3, and CntAttr_{DP4} = 3. The size of the dataset for each data provider $|D_i| = 45, 222$.

Fig. 6 depicts the impact of privacy protection and trust scores on DP₁, DP₂, DP₃, and DP₄'s monetary value. ϵ -differential privacy is enforced with privacy parameters $\epsilon = 0.2, 0.4, 0.6$, and 0.8 and specialization levels $3 \le h \le 19$.

Fig. 6(a) depicts the impact on DP₁, DP₂, DP₃, and DP₄'s monetary value when the threshold is $\epsilon = 0.2$. We observe that DP₄ attains the highest monetary share due to more information utility and its aggregated trust score. When specialization level h increases from 3 to 7 and 11 to 15, DP₁, DP₂, and DP₃ get increases in their monetary shares, while DP₄'s monetary share falls by approximately \$11K, though still achieving a higher share than other data providers. Initially, DP₂ has no monetary share when h = 3, but it increases with the increase in the specialization level h except when h = 19. DP₁, DP₂, and DP₃ is monetary shares become closer to each other when h = 11.

Fig. 6(b) depicts the impact on DP₁, DP₂, DP₃, and DP₄'s monetary value when the threshold is $\epsilon = 0.4$. We observe that DP₄ attains the highest monetary share because of greater information utility and its aggregated trust score. Though DP₁ does not get the highest share, its monetary share becomes closer to DP₄ at h = 11, 15, and 19 with the difference of

approximately \$3K to \$5K. Interestingly, DP₄'s monetary share exhibits nonincreasing monotonicity with the increase in specialization level h, while DP₁'s monetary share increases with the increase in specialization level h except when h = 19. We notice that DP₃ has no monetary share when h = 7 because of a lack of information utility for classification analysis. The trust score does not add any monetary value if a data provider fails to contribute to information utility. The trend on DP₂ and DP₃'s monetary share is not obvious with the increase in h.

Fig. 6(c) depicts the impact on DP₁, DP₂, DP₃, and DP₄'s monetary value when the threshold is $\epsilon = 0.6$. We observe that DP₄ gains the maximum value of monetary share when h = 3 and h = 7, and DP₁ gains the maximum value of monetary share when h = 11 and h = 15, whereas DP₂ gains the maximum value of monetary share when h = 11 and h = 15, whereas DP₂ gains the maximum value of monetary share when h = 19. This is because it has greater information utility in competing with the other data providers at the indicated levels of specialization. We observe that DP₂'s monetary share increases monotonically as the increase in specialization level h, whereas DP₄'s monetary share falls with the increase in specialization level h, except when h = 19.

Fig. 6(d) depicts the impact on DP₁, DP₂, DP₃, and DP₄'s monetary value when the threshold is $\epsilon = 0.8$. We observe that DP₄ achieves the highest monetary share because of greater information utility and its aggregated trust score. We observe that DP₁'s monetary share generally increases as the specialization level *h* increases, whereas DP₄'s monetary share falls with the increase in specialization level *h*, except when h = 11. We notice that when h = 15, all data providers' monetary shares become closer, with a difference of approximately \$4 K.

Case 2: Table IV(b) shows the selection of attributes from each accepted data provider. BA on the integrated data of accepted data providers is 85.4%, using the secure multiple party classifier [21] without disclosing their raw data. We vertically partition the *Adult* dataset into three partitions VP₁, VP₂, and VP₃ for data providers DP₁, DP₂, and DP₄, respectively. Furthermore, we split the dataset into 30 162 and 15 060 records for the training and testing sets, respectively. Since DP₃ has dropped from the list of accepted data providers, DP₄ acquires the position of DP₃. Now, the valuation of each data provider's attribute is \$0.47, \$0.41, and \$0.36, representing ValAttr_{DP1}, ValAttr_{DP2}, and ValAttr_{DP4} by (14), respectively. The attribute count of each data provider is CntAttr_{DP1} = 3, CntAttr_{DP2} = 3,



Fig. 7. Impact of ϵ -differential privacy requirements and trust scores on DP₁, DP₂, and DP₄ monetary value (Case 2). (a) $\epsilon = 0.2$. (b) $\epsilon = 0.4$. (c) $\epsilon = 0.6$. (d) $\epsilon = 0.8$.

and CntAttr_{DP4} = 4. The size of dataset for each data provider $|D_i| = 45222$.

Fig. 7 depicts the impact of privacy protection and trust scores on DP₁, DP₂, and DP₄'s monetary value. ϵ -differential privacy is enforced with privacy parameters $\epsilon = 0.2, 0.4, 0.6, \text{ and } 0.8,$ and specialization levels $3 \le h \le 19$.

Fig. 7(a) depicts the impact on DP₁, DP₂, and DP₄'s monetary value when the threshold is $\epsilon = 0.2$. We observe that DP₄ attains the highest monetary share because of higher information utility and its trust level, except when h = 19. We observe that DP₁'s monetary share increases as the specialization level h increases, except when h = 7, whereas DP₄'s monetary share generally falls with the increase in specialization level h except when h = 15. DP₁ gains the maximum value of approximately \$32K of his monetary share when h = 19.

Fig. 7(b) depicts the impact on DP₁, DP₂, and DP₄'s monetary value when the threshold is $\epsilon = 0.4$. We observe that DP₄ attains the highest monetary share because of higher information utility and its trust level, except when h = 19. The trend on DP₁, DP₂, and DP₄'s monetary share is not obvious with the increase in specialization level h. DP₂ gains the maximum value of approximately \$33K of his monetary share when h = 19.

Fig. 7(c) depicts the impact on DP₁, DP₂, and DP₄'s monetary value when the threshold is $\epsilon = 0.6$. We observe that DP₄ achieves the highest monetary share because of higher information utility and its trust level, except when h = 15. DP₄'s monetary share drops sharply when h increases from 3 to 7 and 11 to 15, while DP₁ and DP₂ have a significant increase in their monetary shares with this increase in h. DP₂ gains the maximum value of approximately \$29K of monetary share when h = 15.

Fig. 7(d) depicts the impact on DP₁, DP₂, and DP₄'s monetary value when the threshold is $\epsilon = 0.8$. We observe that DP₄ gains the maximum value of monetary share when h = 3, 7, and 11, whereas DP₁ gains the maximum value of monetary share when h = 15 and 19. This is because they have more information utility in competing with the other data providers at the indicated levels of specializations. DP₂'s monetary share generally increases as the increase in specialization level h, except when h = 15. DP₁ and DP₄ do not exhibit monotonicity with the increase in h.

VII. CONCLUSION

In this article, we proposed a novel entropy-based trust computation algorithm to verify the correctness of data from untrusted multiple data providers who own overlapping attributes over the same set of records. We achieved three main benefits in delegating the verification role to the semitrusted CSP. First, our method ensured that the CSP cannot derive customers' private data from the information collected during the verification process. Second, the overhead of computation on the CS was also reduced because only an encrypted information gain message and its keyed hash were exchanged between a data provider and the CS, instead of exchanging encrypted individual data records during the verification process. Third, it also reduced the burden on data consumers to determine, which data providers can serve their demands on requested attributes and what are their attained trust scores. Furthermore, we evaluated the robustness of our approach when a data provider employed the machine learning method for imputation of missing values on its data. There was no significant difference in the perspective to the performance of the imputation method. It is conditional to what proportion of data is missing and whether the data contains repeated patterns. If the prediction of a missing data happens to be as precise data, then it will be considered as true data. We incorporated the VCG auction mechanism to determine the pricing on data providers' attributes. It maximized the total valuation obtained by data providers, since there was no incentive to lie or deviate from truthful reporting for a data provider. From the perspective of privacy protection, the accepted data providers as a result of trust computation set up their joint privacy requirements for the data mashup. During the data mashup process, every data provider competed with the other participating data providers to produce more data utility. It was evident from the experiments that an accepted data provider whose data attributes result in more information gain, and whose trust level is higher than the other competitors, can get a significantly larger share of the monetary value.

REFERENCES

- "Data Partners," Seventh Point, 2014. Accessed: Jun. 11, 2015. [Online]. Available: http://www.seventhpoint.com/whitepaper/data-partners/
- [2] O. A. Wahab, J. Bentahar, H. Otrok, and A. Mourad, "A survey on trust and reputation models for web services: Single, composite, and communities," *Decis. Support Syst.*, vol. 74, pp. 121–134, 2015.
- [3] R. Agrawal, A. Evfimievski, and R. Srikant, "Information sharing across private databases," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2003, pp. 86–97.
- [4] D. Alhadidi, N. Mohammed, B. C. M. Fung, and M. Debbabi, "Secure distributed framework for achieving ε-differential privacy," in *Proc. 12th Int. Symp. Privacy Enhancing Technol.*, 2012, pp. 120–139.
- [5] R. R. Andridge and R. J. A. Little, "A review of hot deck imputation for survey non-response," *Int. Statist. Rev.*, vol. 78, no. 1, pp. 40–64, 2010.

- [6] M. Anisetti, C. A. Ardagna, and E. Damiani, "A certification-based trust model for autonomic cloud computing systems," in Proc. Int. Conf. Cloud Autonomic Comput., 2014, pp. 212-219.
- [7] Y. Aumann and Y. Lindell, "Security against covert adversaries: Efficient protocols for realistic adversaries," J. Cryptol., vol. 23, no. 2, pp. 281-343, 2010.
- [8] O. Benjelloun, A. Das Sarma, A. Halevy, M. Theobald, and J. Widom, "Databases with uncertainty and lineage," VLDB J., vol. 17, no. 2, pp. 243-264, 2008
- [9] L. Beretta and A. Santaniello, "Nearest neighbor imputation algorithms: A critical evaluation," BMC Med. Inform. Decis. Making, vol. 16, no. 3, 2016, Art. no. 74.
- [10] B. V. D. Berg and E. Keymolen, "Regulating security on the Internet: Control versus trust," Int. Rev. Law, Comput. Technol., vol. 31, no. 2, pp. 188-205, 2017.
- [11] E. Bertino and H.-S. Lim, "Assuring data trustworthiness: Concepts and research challenges," in Proc. 7th VLDB Conf. Secure Data Manage., 2010, рр. 1–12.
- [12] E. Bertino, L. Martino, F. Paci, and A. Squicciarini, Security for Web Services and Service-Oriented Architectures, 1st ed. Berlin, Germany: Springer, 2009.
- [13] E. Bertino and R. Sandhu, "Database security-Concepts, approaches, and challenges," IEEE Trans. Dependable Secure Comput., vol. 2, no. 1, pp. 2-19, Jan. 2005.
- [14] D. Boneh and V. Shoup, A Graduate Course in Applied Cryptography. Stanford, CA, USA: Stanford Univ., 2017.
- [15] V. Chang, Y.-H. Kuo, and M. Ramachandran, "Cloud computing adoption framework: A security framework for business clouds," Future Gener. Comput. Syst., vol. 57, pp. 24-41, 2016.
- [16] C. Dai, D. Lin, E. Bertino, and M. Kantarcioglu, "An approach to evaluate data trustworthiness based on data provenance," in Proc. Workshop Secure Data Manage., 2008, pp. 82-98.
- [17] E. De Cristofaro, P. Gasti, and G. Tsudik, "Fast and private computation of cardinality of set intersection and union," in Proc. 11th Int. Conf. Cryptol. Netw. Secur., 2012, pp. 218-231.
- [18] T. Dierks and E. Rescorla, "The transport layer security (TLS) protocol version 1.2," RFC 5246, 2008.
- C. Dong, L. Chen, J. Camenisch, and G. Russello, "Fair private set [19] intersection with a semi-trusted arbiter," in Proc. 27th Int. Conf. Data Appl. Secur. Privacy, 2013, pp. 128-144.
- [20] C. Dong, L. Chen, and Z. Wen, "When private set intersection meets big data: An efficient and scalable protocol," in Proc. ACM SIGSAC Conf. Comput. Commun. Secur., 2013, pp. 789-800.
- [21] W. Du and Z. Zhan, "Building decision tree classifier on private data," in Proc. IEEE Int. Conf. Privacy, Secur. Data Mining, 2002, vol. 14, pp. 1-8.
- [22] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating noise to sensitivity in private data analysis," in Proc. 3rd Conf. Theory Cryptography, 2006, pp. 265-284.
- [23] D. Easley and J. Kleinberg, Networks, Crowds, and Markets: Reasoning About a Highly Connected World. Cambridge, U.K.: Cambridge Univ. Press, 2010.
- [24] The 2018 State of Data Management: A Public Sector Benchmark Report, Experian Information Solutions, Inc., Costa Mesa, CA, USA, 2018.
- [25] C. Feijo, J. L. Gmez-Barroso, and P. Voigt, "Exploring the economic value of personal information from firms' financial statements," Int. J. Inf. Manage., vol. 34, no. 2, pp. 248-256, 2014.
- [26] M. J. Freedman, K. Nissim, and B. Pinkas, "Efficient private matching and set intersection," in Proc. Int. Conf. Theory Appl. Cryptographic Techn., 2004, pp. 1-19.
- [27] J. Freudiger, S. Rane, A. E. Brito, and E. Uzun, "Privacy preserving data quality assessment for high-fidelity data sharing," in Proc. ACM Workshop Inf. Sharing Collaborative Secur., 2014, pp. 21-29.
- [28] A. Friedman and A. Schuster, "Data mining with differential privacy," in Proc. 16th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 2010, p. 493-502.
- [29] B. C. M. Fung, T. Trojer, P. C. Hung, L. Xiong, K. Al-Hussaeni, and R. Dssouli, "Service-oriented architecture for high-dimensional private data mashup," IEEE Trans. Services Comput., vol. 5, no. 3, pp. 373-386, Jul.-Sep. 2012.
- [30] B. C. M. Fung, K. Wang, R. Chen, and P. S. Yu, "Privacy-preserving data publishing: A survey of recent developments," ACM Comput. Surv., vol. 42, no. 4, 2010, Art. no. 14.
- [31] B. C. M. Fung, K. Wang, and P. S. Yu, "Anonymizing classification data for privacy preservation," IEEE Trans. Knowl. Data Eng., vol. 19, no. 5, pp. 711-725, 2007.

- [32] S. R. Ganta, S. P. Kasiviswanathan, and A. Smith, "Composition attacks and auxiliary information in data privacy," in Proc. 14th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 2008, pp. 265-273.
- [33] Y. Huang, D. Evans, and J. Katz, "Private set intersection: Are garbled circuits better than custom protocols?" in Proc. 19th Netw. Distrib. Syst. Secur. Symp., 2012. [Online]. Available: https://dblp.org/db/conf/ndss/ ndss2012.html
- [34] L. A. Hunt, Missing Data Imputation and Its Effect on the Accuracy of Classification. Berlin, Germany: Springer, 2017.
- [35] Y. Ishai, J. Kilian, K. Nissim, and E. Petrank, "Extending oblivious transfers efficiently," in Proc. Annu. Int. Cryptol. Conf., 2003, pp. 145-161.
- [36] W. Jiang and C. Clifton, "A secure distributed framework for achieving K-anonymity," VLDB J., vol. 15, no. 4, pp. 316-333, 2006.
- [37] P. Jurczyk and L. Xiong, "Distributed anonymization: Achieving privacy for both data subjects and data providers," in Proc. 23rd Annu. IFIP WG 11.3 Working Conf. Data Appl. Secur., 2009, pp. 191-207.
- [38] A. Jsang, R. Ismail, and C. Boyd, "A survey of trust and reputation systems for online service provision," Decis. Support Syst., vol. 43, no. 2, pp. 618-644, 2007.
- [39] S. Kamara, P. Mohassel, M. Raykova, and S. Sadeghian, "Scaling private set intersection to billion-element sets," in Proc. Financial Cryptography Data Secur., 2014, pp. 195-215.
- [40] S. Kamara, P. Mohassel, and B. Riva, "Salus: A system for server-aided secure function evaluation," in Proc. ACM Conf. Comput. Commun. Secur., 2012, pp. 797-808.
- [41] S. Karabati and Z. B. Yalcin, "An auction mechanism for pricing and capacity allocation with multiple products," Prod. Oper. Manage., vol. 23, no. 1, pp. 81-94, 2014.
- [42] V. Kher and Y. Kim, "Securing distributed storage: Challenges, techniques, and systems," in Proc. ACM Workshop Storage Secur. Survivability, 2005, pp. 9-25.
- R. H. Khokhar, B. C. M. Fung, F. Iqbal, D. Alhadidi, and J. Benta-[43] har, "Privacy-preserving data mashup model for trading person-specific information," Electron. Commerce Res. Appl., vol. 17, pp. 19-37, 2016.
- [44] D. Kifer, "Attacks on privacy and deFinetti's theorem," in Proc. ACM SIGMOD Int. Conf. Manage. Data, 2009, pp. 127-138.
- [45] D. J. Kim, D. L. Ferrin, and H. R. Rao, "A trust-based consumer decisionmaking model in electronic commerce: The role of trust, perceived risk, and their antecedents," Decis. Support Syst., vol. 44, no. 2, pp. 544-564, 2008.
- [46] V. Kolesnikov, R. Kumaresan, M. Rosulek, and N. Trieu, "Efficient batched oblivious PRF with applications to private set intersection," in Proc. ACM SIGSAC Conf. Comput. Commun. Secur., 2016, pp. 818-829.
- N. Li, T. Li, and S. Venkatasubramanian, "t-closeness: Privacy beyond [47] k-Anonymity and l-Diversity," in Proc. 23rd IEEE Int. Conf. Data Eng., 2007, pp. 106-115.
- [48] X. Li, R. Ding, X. Liu, X. Liu, E. Zhu, and Y. Zhong, "A dynamic pricing reverse auction-based resource allocation mechanism in cloud workflow systems," Sci. Program., vol. 2016, pp. 1-13, 2016.
- [49] H.-S. Lim, G. Ghinita, E. Bertino, and M. Kantarcioglu, "A game-theoretic approach for high-assurance of data trustworthiness in sensor networks." in Proc. 28th IEEE Int. Conf. Data Eng., 2012, pp. 1192-1203.
- [50] H.-S. Lim, Y.-S. Moon, and E. Bertino, "Provenance-based trustworthiness assessment in sensor networks," in Proc. 7th Int. Workshop Data Manage. Sens. Netw., 2010, pp. 2-7.
- A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkitasubramaniam, " ℓ -diversity: Privacy beyond *k*-anonymity," *ACM Trans. Knowl. Discovery* [51] Data, vol. 1, no. 1, 2007, doi: 10.1145/1217299.1217302.
- [52] F. McSherry and K. Talwar, "Mechanism design via differential privacy," in Proc IEEE . 48th Annu. Symp. Found. Comput. Sci., 2007, pp. 94-103.
- [53] N. Mohammed, D. Alhadidi, B. C. M. Fung, and M. Debbabi, "Secure two-party differentially private data release for vertically partitioned data," IEEE Trans. Dependable Secure Comput., vol. 11, no. 1, pp. 59-71, Jan./Feb. 2014.
- [54] N. Mohammed, B. C. M. Fung, and M. Debbabi, "Anonymity meets game theory: secure data integration with malicious participants," VLDB J., vol. 20, no. 4, pp. 567-588, 2011.
- [55] N. Mohammed, B. C. M. Fung, P. C. K. Hung, and C.-K. Lee, "Centralized and distributed anonymization for high-dimensional healthcare data," ACM Trans. Knowl. Discovery Data, vol. 4, no. 4, 2010, Art. no. 18.
- R. Nget, Y. Cao, and M. Yoshikawa, "How to balance privacy and money [56] through pricing mechanism in personal data market," in Proc. SIGIR Workshop eCommerce, vol. 2311, 2017. [Online] Available: http://ceurws.org/Vol-2311/

- [57] T. H. Noor, Q. Z. Sheng, L. Yao, S. Dustdar, and A. H. Ngu, "CloudArmor: Supporting reputation-based trust management for cloud services," *IEEE Trans. Parallel Distrib. Syst.*, vol. 27, no. 2, pp. 367–380, Feb. 2016.
- [58] OECD, "Exploring the economics of personal data: A survey of methodologies for measuring monetary value," OECD Digit. Economy Papers 220, 2013.
- [59] B. Pinkas, T. Schneider, and M. Zohner, "Scalable private set intersection based on OT extension," ACM Trans. Privacy Secur., vol. 21, no. 2, 2018, Art. no. 7.
- [60] B. Pinkas, T. Schneider, G. Segev, and M. Zohner, "Phasing: Private set intersection using permutation-based hashing," in *Proc. 24th USENIX Conf. Secur. Symp.*, 2015, pp. 515–530.
- [61] B. Pinkas, T. Schneider, and M. Zohner, "Faster private set intersection based on OT extension," in *Proc. 23rd USENIX Conf. Secur. Symp.*, 2014, pp. 797–812.
- [62] Y. Pu and J. Grossklags, "Valuating friends' privacy: Does anonymity of sharing personal data matter?" in *Proc. 13th Symp. Usable Privacy Secur.*, 2017, pp. 339–355.
- [63] J. R. Quinlan, "Induction of decision trees," *Mach. Learn.*, vol. 1, no. 1, pp. 81–106, 1986.
- [64] D. Quintero et al., IBM Software Defined Environment, 1st ed. Armonk, NY, USA: IBM Redbooks, 2015.
- [65] E. Rescorla, "The Transport Layer Security (TLS) Protocol Version 1.3," RFC 8446, 2018.
- [66] P. Samarati, "Protecting Respondents' Identities in Microdata Release," *IEEE Trans. Knowl. Data Eng.*, vol. 13, no. 6, pp. 1010–1027, 2001.
- [67] R. Shaikh and M. Sasikumar, "Trust model for measuring security strength of cloud computing service," *Procedia Comput. Sci.*, vol. 45, pp. 380–389, 2015.
- [68] C. E. Shannon, "The Mathematical Theory of Communication," *Bell Syst. Tech. J.*, vol. 27, no. 3, pp. 379–423, 1948.
- [69] L.-A. Tang, X. Yu, S. Kim, J. Han, C.-C. Hung, and W.-C. Peng, "Tru-Alarm: Trustworthiness analysis of sensor networks in cyber-physical systems," in *Proc. IEEE Int. Conf. Data Mining*, 2010, pp. 1079–1084.
- [70] M. Tang, Y. Xu, J. Liu, Z. Zheng, and X. Liu, "Combining global and local trust for service recommendation," in *Proc. IEEE Int. Conf. Web Services*, 2014, pp. 305–312.
- [71] R. C.-W. Wong, A. W.-C. Fu, K. Wang, and J. Pei, "Minimality attack in privacy preserving data publishing," in *Proc. 33rd Int. Conf. Very Large Data Bases*, 2007, pp. 543–554.
- [72] Q. Wu, M. Zhou, Q. Zhu, and Y. Xia, "VCG Auction-Based dynamic pricing for multigranularity service composition," *IEEE Trans. Autom. Sci. Eng.*, vol. 15, no. 2, pp. 796–805, Apr. 2018.
- [73] X. Xiao, G. Bender, M. Hay, and J. Gehrke, "iReduct: Differential privacy with reduced relative errors," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2011, pp. 229–240.
- [74] Y. Xiao, L. Xiong, L. Fan, S. Goryczka, and H. Li, "DPCube: Differentially private histogram release through multidimensional partitioning," *Trans. Data Privacy*, vol. 7, no. 3, pp. 195–222, 2014.
- [75] X. Yang, X. Luo, X. A. Wang, and S. Zhang, "Improved outsourced private set intersection protocol based on polynomial interpolation," *Concurrency Comput.: Pract. Exp.*, vol. 30, no. 1, 2017, Art. no. e4329.
- [76] L. Zhang and P. N. Suganthan, "Random forests with ensemble of feature spaces," *Pattern Recognit.*, vol. 47, no. 10, pp. 3429–3437, 2014.
- [77] Z. Zhang, "Missing data imputation: Focusing on single imputation," Ann. Transl. Med., vol. 4, no. 1, 2016, doi: 10.3978/j.issn.2305-5839.2015.12.38.



Rashid Hussain Khokhar received the master's degree in information systems security in 2013 from the Concordia Institute for Information Systems Engineering, Concordia University, Montreal, QC, Canada, where he is currently working toward the Ph.D. degree in information and systems engineering.

He served as a Reviewer for prestigious conferences, including the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, the IEEE International Conference on Data Min-

ing, the International Joint Conference on Artificial Intelligence, and the Pacific-Asia Conference on Knowledge Discovery and Data Mining. His current research interests include services computing, privacy-preserving data publishing, data mining and machine learning, data trustworthiness, information security, and risk-benefit analysis.

Mr. Khokhar received an International Tuition Fee Remission Award for his Ph.D. studies.



Farkhund Iqbal received the master's and Ph.D. degrees in computer science from Concordia University, Montreal, QC, Canada, in 2005 and 2011, respectively.

He is currently an Associate Professor and Director of the Advanced Cyber Forensics Research Laboratory, College of Technological Innovation, Zayed University, Abu Dhabi, United Arab Emirates. He has authored or coauthored more than 100 papers in high-ranked journals and conferences. He is an Affiliate Professor with the School of Information

Studies, McGill University, Montreal and an Adjunct Professor with the Faculty of Business and IT, Ontario Tech University, Oshawa, ON, Canada. He is the recipient of several prestigious awards and research grants. He has served as a Chair and Technical Program Committee Member of several IEEE/ACM conferences and is the Reviewer of high-rank journals. His research interests include machine learning and big data techniques for problem solving in health care, cybersecurity, and cybercrime investigation in the smart and safe city domain.



Benjamin C. M. Fung (Senior Member, IEEE) received the Ph.D. degree in computing science from Simon Fraser University, Burnaby, BC, Canada, in 2007.

He is currently a Canada Research Chair in Data Mining for Cybersecurity, an Associate Professor with the School of Information Studies, an Associate Member with the School of Computer Science, McGill University, Montreal, QC, Canada, and a Co-Curator of Cybersecurity in the World Economic Forum. He has authored or coauthored more than 120

refereed publications that span the research forums of data mining, privacy protection, cybersecurity, services computing, and building engineering. His data mining works in crime investigation and authorship analysis have been reported by media worldwide.

Dr. Fung is a licensed Professional Engineer of Software Engineering in Ontario, Canada.



Jamal Bentahar (Member, IEEE) received the bachelor's degree in software engineering from the National Institute of Statistics and Applied Economics, Rabat, Morocco, in 1998, the M.Sc. degree in software engineering from Mohamed V University, Rabat, in 2001, and the Ph.D. degree in computer science and software engineering from Laval University, Quebec City, QC, Canada, in 2005.

He is currently a Professor with the Concordia Institute for Information Systems Engineering, Faculty of Engineering and Computer Science, Concordia

University, Montreal, QC, Canada. From 2005 to 2006, he was a Postdoctoral Fellow with Laval University and then a Natural Sciences and Engineering Research Council Postdoctoral Fellow with Simon Fraser University, Burnaby, BC, Canada. His research interests include computational logics, model checking, multiagent systems, service and cloud computing, game theory, and software engineering. He is an NSERC Co-Chair for Discovery Grant for Computer Science (2016–2018).