

A privacy-preserving statistics marketplace using local differential privacy and blockchain: an application to smart-grid measurements sharing

Nikos Fotiou, Iakovos Pittaras, Vasilios A. Siris, George C. Polyzos

*Mobile Multimedia Laboratory
Department of Informatics, School of Information Sciences and Technology
Athens University of Economics and Business, Greece*

Priit Anton

*Guardtime
Tallinn, Estonia*

Abstract

Service providers usually require detailed statistics in order to improve their services. On the other hand, privacy concerns are intensifying and sensitive data is protected by legislation, such as GDPR. In this paper we present the design, implementation, and evaluation of a marketplace that allows “data consumers” to buy information from “data providers,” which can then be used for generating meaningful statistics. Additionally, our system enables “system operators” that can select which data providers are allowed to provide data, based on filtering criteria specified by the data consumer. We leverage local differential privacy to protect the data provider’s privacy against data consumers, as well as against system operators, and we build a blockchain-based solution for ensuring fair exchange, and immutable data logs. Our design targets use cases that involve hundreds or even thousands of data providers. We prove the feasibility of our approach through a proof-of concept implementation of a measurement sharing application for smart-grid systems.

Keywords: Auditability, Ethereum, Fair exchange, Smart contracts

Email addresses: {fotiou,pittaras,vsiris,polyzos}@aueb.gr (Nikos Fotiou, Iakovos Pittaras, Vasilios A. Siris, George C. Polyzos), priit.anton@guardtime.com (Priit Anton)

1 **1. Introduction**

2 We are living in a globalized, cyber connected society where users have a
3 plethora of choices. In this highly competitive environment, service providers
4 try to offer as much personalized and consumer-tailored services as possible. In
5 order to achieve their goal, they seek access to user profiling information. How-
6 ever, increased privacy concerns, as well as legislation, such as EU’s General
7 Data Protection Regulation (GDPR), have made the collection of such informa-
8 tion a thorny challenge. Of course this comes as no surprise, since such activities
9 not only jeopardize users’ privacy, but also, as we have recently witnessed, the
10 collected information can be used for manipulating users’ choices [1]. Hence,
11 the research question “how can sensitive data be securely shared?” still remains
12 open.

13 Traditionally, the collection of sensitive information has been protected using
14 anonymization techniques. However, large-scale privacy breaches from suppos-
15 edly anonymized datasets, such as those involving AOL [2] and Netflix [3], have
16 questioned the ability of those techniques to effectively protect user privacy [4].
17 A more promising approach is the use of *differential privacy* [5]. The goal of
18 differential privacy is to allow the extraction of statistics about a population
19 of users without revealing any information about specific individuals. This is
20 achieved with the addition of some “noise” to the statistics extraction process.
21 The added noise guarantees that the contribution of a single individual to the
22 calculated statistics is not “significant.” In other words, differential privacy
23 guarantees that the output of the statistics calculation process is only slightly
24 impacted by the contributions of a single individual, hence meaningful informa-
25 tion about these contributions cannot be extracted.

26 Most related systems consider the so-called “centralized” differential privacy
27 approach, where a trusted 3rd party collects user data adds the appropriate
28 amount of noise, and releases privacy-preserving statistics. Although this ap-
29 proach produces accurate results even with few samples, the introduction of a

30 trusted entity is a significant security and privacy risk. Furthermore, there can
31 be cases where regulations or laws prohibit such 3rd parties.

32 In this paper, we propose a privacy-preserving solution that allows a “data
33 consumer” to extract meaningful statistics from sensitive data protected using
34 “local differential privacy.” Local differential privacy enables “data providers” to
35 add noise to their (sensitive) data by themselves, and share them with untrusted
36 3rd party data consumers without jeopardizing their privacy. Additionally, our
37 solution considers an intermediate entity, referred to as the “system operator,”
38 that coordinates the whole process and applies filtering rules. With our ap-
39 proach, and as opposed to the state of the art, data providers are protected
40 even against “curious” system operators. This is achieved by having the data
41 providers send their (noisy) data directly to the data consumers.

42 Nevertheless, our approach leads inevitably to some tussles. For instance,
43 a data consumer may be tempted to not pay the required fee. Similarly, data
44 providers may indicate their interest to participate in a data collection pro-
45 cess, receive the corresponding fee, but refuse to provide the actual data. Fi-
46 nally, a system operator may incorrectly filter out the responses of certain data
47 providers. In order to resolve these tussles we rely on the blockchain technol-
48 ogy. In particular, we leverage Ethereum smart contracts to provide a privacy-
49 preserving immutable log of operations that can be used for dispute resolution,
50 as well as to provide “fair exchange” of data and service fees. The contributions
51 of our work presented in this paper are:

- 52 • We allow a “data consumer” to specify a minimum number of “data
53 providers” it wishes to receive statistics from: if fewer data providers
54 participate in an exchange then the data consumer does not have to pay
55 for a service fee but at the same time it cannot generate any statistics.
- 56 • We enable “system operators” to filter out data consumers based on some
57 specified criteria, without having access to their (noisy) data.
- 58 • We enable “fair exchange” between a data consumer and many (hundreds
59 or even thousands) data providers, i.e., with our solution a data consumer

60 can access the (noisy) data of all providers only after paying the required
61 service fee.

- 62 • Even though our solution involves many stakeholders, the cost for using
63 the blockchain technology is minimum since the smart contract records
64 only the minimum information necessary for auditing.

65 In order to put our solution in context, we consider a smart-grid measure-
66 ments sharing use case.

67 With the smart meter penetration rate reaching 42,5% by 2020 and 83,87%
68 by 2024 in the EU-28 [6], metering data can be a valuable source of informa-
69 tion to various service providers. Although many *Distribution System Operators*
70 (DSOs) provide data hubs for metering data, the vast majority of them prohibit
71 3rd party service providers (e.g., telecom operators, service partners, techni-
72 cal integrators, and other 3rd parties that may provide added-value services to
73 smart meter owners) from accessing the stored data, mainly due to privacy con-
74 cerns. Using our approach, smart meters can add noise to their measurements
75 and directly share them with 3rd parties, alleviating the burden of maintaining
76 sensitive data from DSOs.

77 The rest of the paper is organized as follows. In Section 2 we provide back-
78 ground information, as well as related work in this area. In Section 3 we give
79 an overview of the proposed system, and we detail its design in Section 4. We
80 present the evaluation of our system in Section 5. Finally, we discuss some prop-
81 erties of the system in Section 6, and we present our conclusions in Section 7.

82 **2. Background and Related work**

83 *2.1. Smart contracts and fair exchange*

84 Smart contracts are decentralized applications that are deterministically ex-
85 ecuted in a blockchain (e.g., Ethereum). Smart contracts are widely used for
86 implementing *escrow* services for exchanging digital goods in a fair way [7]. In
87 particular, such an escrow smart contract allows “buyers” to deposit digital cur-
88 rency, (a portion of) which is transferred to a “seller” when a proof of digital

Are you a member of the government?

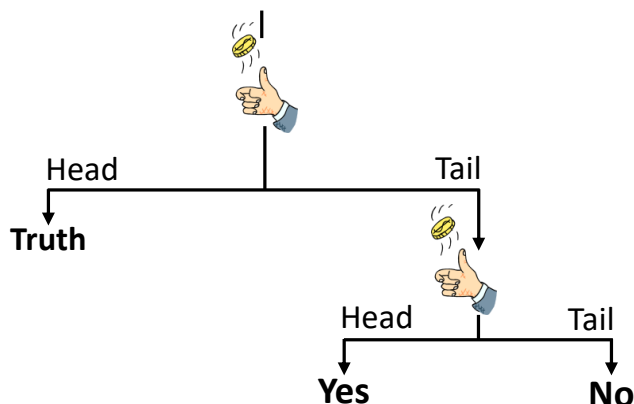


Figure 1: An example of randomized response. A respondent says the truth only if the flip of the coin comes up head, otherwise her response is determined by a second flip of the coin.

89 good exchange is presented to the contract. From a high-level perspective this
90 exchange is implemented as follows: a buyer wishes to buy an information item
91 i ; the escrow smart contract is configured with the hash $h(i)$ of that item; the
92 buyer deposits the appropriate amount of digital currency in the smart contract;
93 then, the seller reveals to the smart contract an item i' : the smart contract ver-
94 ifies that $h(i) = h(i')$, and if the verification is successful (which means $i' = i$)
95 it transfers (a portion of) the deposited digital currency to the seller. Since the
96 information recorded in a smart contract is public, the buyer can retrieve i' .

97 2.2. Local differential privacy

98 The main building block of local differential privacy mechanisms is an old
99 survey technique known as “Randomized Responses” [8]. This technique in-
100 volves a (sensitive) question that can be answered with a “Yes” or “No” (e.g.,
101 “Do you work for the government?”). A respondent flips a fair coin in secret,
102 if the coin comes up heads she responds honestly, otherwise she flips the coin
103 in secret again; this time if the coin comes up heads she responds with “Yes”
104 and if it comes up tails she responds with “No”. It can be seen that with this
105 process, illustrated in Figure 1, 75% of the time the response will correspond

106 to the respondent’s true answer. Moreover, and provided that she is not asked
107 again the same question, the respondent’s privacy is preserved since 50% of the
108 time her response is generated at random. After a number of responses has
109 been collected, the probability of a “Yes” response can be accurately estimated
110 by calculating $\text{minimum}(0, 2 * (Y - 0.25))$, where Y is the proportion of “Yes”
111 in the submitted responses [9].

112 RAPPOR [9], which has been developed by Google and is included with the
113 Chrome browser, extends “Randomized Responses” surveys to allow for multiple
114 possible answers. In its basic form (“Basic one-time RAPPOR,” Section 2.1
115 of [9]), which is considered in this paper, RAPPOR involves a question that can
116 be answered with a pre-defined number of choices (e.g., “Which of the following
117 operating systems do you use: a) Linux b) Windows c) Mac OS d) Other ”).
118 A respondent constructs a bit vector of size equal to the number of choices
119 (i.e., 4 in our example). The first position of the bit vector corresponds to the
120 first choice, the second position to the second choice, and so forth. Then for
121 each choice, she executes the randomized response algorithm (i.e, first she will
122 respond to “do you use Linux?”, then to “do you use Windows?” and so forth).
123 If the output of the algorithm is “Yes”, she fills the corresponding position of
124 the vector with “1”, otherwise with “0”. After a number of bit vectors has been
125 collected, the probability of the i^{th} choice is again estimated by calculating
126 $\text{minimum}(0, 2 * (Y - 0.25))$, where Y is now the proportion of ones in the i^{th}
127 position of the collected bit vectors (e.g., for the calculation of the probability
128 of the choice “a) Linux” Y is set to the proportion of ones in the first position
129 of all bit vectors).

130 2.3. Related work

131 Differential privacy has been considered by many research efforts including
132 recommendation systems [10], data mining [11], crowd-sourcing [12], tools for
133 performing network measurements [13], intelligent transportation systems [14],
134 sensor network stream processing [15], and many others. These solutions have
135 two key differences compared to our approach, firstly the differential privacy

136 mechanisms are applied by a trusted entity; this entity collects data from users,
137 performs some computations, and extracts some privacy preserving statistics
138 by adding “noise”; secondly, they assume that the entity that generates the
139 statistics has also some sort of business relationship with the users, hence it is
140 capable of “filtering” users based on pre-defined criteria (e.g., they can calculate
141 the average of some metric only for users satisfying a location requirement).
142 Unlike the above, our solution applies differential privacy at the data source,
143 hence offers better privacy (however, this higher privacy comes at the cost of
144 requiring more data inputs to achieve similar accuracy to that of centralized
145 approaches). Secondly, our solution differentiates the entities that calculate a
146 statistic from those that perform user filtering. This has the advantage that the
147 entity responsible for filtering does not learn the calculated statistic. In order
148 to illustrate the importance of this property consider the case of a university
149 ranking authority; assume this authority wants to calculate how students rank
150 their universities: whenever a user submits a response, the university can verify
151 whether or not this user is a student, but it does not have access to the response
152 submitted by the user nor to the calculated rank.

153 Recently, a number of research efforts have leveraged blockchain technology
154 to implement “fair-exchange” of data, focusing at the same time on privacy
155 preservation. Dimitriou and Mohammed [16] are using Bitcoin to perform fair
156 exchange of smart-grid measurements. Nevertheless, with their solution, once
157 the payment is made the data consumer has access to the actual user data,
158 hence user privacy is not preserved. In contrast, with our solution user privacy is
159 always protected. Another notable difference is that the solution in [16] performs
160 an “one-to-one” fair exchange. Our approach supports a more complex model
161 where many users encrypt their data: once a sufficient number of users have
162 provided data, the data consumer makes the payment, and the (noisy) data
163 of all users is revealed. Duan et al. [17] also use a blockchain-based approach
164 where users stored their encrypted data in an Ethereum smart contract; then,
165 a trusted entity decrypts the data, adds some noise (i.e., centralized differential
166 privacy is used), and performs a fair exchange with a data consumer. Unlike this

167 approach, with our solution, users do not record their data in the blockchain
168 (not even encrypted). This design choice is driven by the fact that our solution
169 considers hundreds (or even thousands) of users: if users were recording their
170 data in the smart contract the total transaction cost would be prohibitive.

171 A number of research efforts propose computations over sensitive informa-
172 tion using cryptographic techniques, such as secure multi-party computation,
173 or fully-homomorphic encryption (e.g., [18]). An advantage of these solutions,
174 compared to our approach, is that they achieve accurate statistics. Of course,
175 this is a double-edged sword: accurate statistics introduce privacy risks when
176 the number of samples is small, or in cases where an attacker is able to monitor
177 how the computed statistic “evolves” with the addition (or removal) of samples,
178 which is the threat model of differential privacy. Moreover, these solutions have
179 high complexity and impose significant computational overhead, hence they are
180 not suitable for all use cases.

181 Traditional differential privacy mechanisms are based on adding noise sam-
182 pled from a Laplace distribution to the aggregated data. In order for this noise
183 to be added in a distributed manner, and hence avoid the need for a trusted
184 aggregator, data providers should be able to communicate with each other. For
185 example, Acs and Castelluccia [19] build a differential privacy solution for smart
186 meter aggregate statistics (i.e., the same use case as in our paper). In order to
187 avoid trusted aggregators, they organize smart meters in clusters and require
188 the meters of each cluster to communicate with each other. Our solution does
189 not require any communication among data providers, neither does it require
190 that data providers know any information about the set of providers respond-
191 ing to a query (e.g., the total number of providers that send measurements).
192 Similarly, a number of related works add noise locally, to a series of data, and
193 the amount of noise is optimized based on the previous records (e.g. see [20]).
194 Our solution assumes a different use case: data providers provide a single data
195 point (e.g., consumption in a particular day). Instead of achieving local differ-
196 ential privacy using the Laplace distribution, our system uses the much simpler
197 model of randomized responses. To this end, we build on RAPPOR, but other

198 randomized response systems can be used instead (e.g., a more recent system
199 presented in [21]).

200 Gai et al. [22] use differential privacy to build a privacy preserving blockchain-
201 based architecture for Industrial IoT. The goal of this architecture is to provide a
202 privacy-preserving task allocation to “edge devices.” The proposed architecture
203 relies on a centralized entity referred to as the “Optimization Server” which is
204 responsible for assigning tasks, collecting data, and adding “noise” to data using
205 differential privacy techniques. In contrast, our approach uses local differential
206 privacy, hence it does not need such a trusted centralized entity.

207 Fioreto et al. [23] apply differential privacy to obfuscate power line param-
208 eters without preventing however grid optimization algorithms. The proposed
209 solution is very specific to the particular problem (power line parameters obfus-
210 cation) and it assumes a globally accessible model of the network to solve the
211 problem. Our solution is simpler since it does not require any “global knowl-
212 edge” from the data providers. Nevertheless, our solution cannot be used for
213 the purposes of the problem discussed in [23]: the obfuscated response of a
214 data consumer cannot be used as an input to an algorithm that requires high
215 precision data (e.g., an optimization algorithm).

216 A number of works assume that smart meters cannot be trusted/modified
217 and for this reason they use other devices operating in parallel to achieve local
218 differential privacy. For example Zhao et al. [24] use batteries and through their
219 charging and discharging functions they hide the energy usage patterns of the
220 household. Our paper assumes that a smart meter is programmable and hence
221 it can be modified to execute our local differential privacy algorithm.

222 Related work in this area investigates the integration of differential privacy
223 into blockchain protocols. For example Hassan et al. [25] discuss how differential
224 privacy can be integrated into blockchain building blocks. But this is completely
225 orthogonal to our approach; we do not propose any modification to blockchain
226 technology and we do not even store data in the blockchain (only hashes).

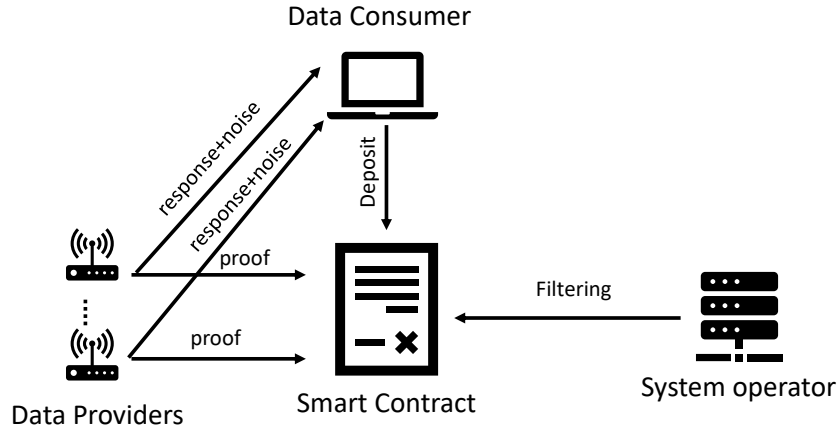


Figure 2: System overview

227 **3. System overview and goals**

228 Our system is composed of the following entities (Figure 2): A *data con-*
 229 *sumer* interested in extracting statistics about a population that meets certain
 230 criteria. Many *data providers* that can provide input used for generating the de-
 231 sired statistics. A *system operator* that is responsible for “filtering” the provided
 232 responses, i.e., for determining if the corresponding data provider meets the cri-
 233 teria specified by the data consumer. From a high-level perspective these entities
 234 interact as follows: A data consumer constructs a “query” which is answered by
 235 a data provider using local differential privacy. Such a query can be regarded
 236 as a vector $Q = \{c_1, c_2, \dots, c_n\}$, where c_1, c_2, \dots, c_n are possible choices. Addi-
 237 tionally, the data consumer specifies “criteria” that define which data providers
 238 are allowed to respond. Data providers “express their interest” to respond to a
 239 query. The system operator indicates which of the data providers that expressed
 240 interest meet the defined criteria. Then, those providers send their responses
 241 directly to the data consumer. Eventually, the data consumer extracts the de-
 242 sired statistics, compensates the system operator, which in turn compensates
 243 the appropriate data providers. The extracted statistics indicate the percentage
 244 of data providers that responded positively to each choice $Q[c_x] \forall x \in [1, n]$.

245 As we discuss in the following section, to achieve a high accuracy of the

246 extracted statistic, the number of responses should be of the order of hundreds
247 or even thousands. Our system has the following properties:

- 248 • Nobody, not even the system operator can determine the exact response
249 of a data provider.
- 250 • A data consumer cannot extract any statistics before paying for the pro-
251 vided services.
- 252 • A data consumer is required to pay only if a pre-agreed number of re-
253 sponses has been collected.

254 Our system achieves these goals through a blockchain-based solution. In
255 particular, we design an Ethereum smart contract that can be used as an “im-
256 mutable log” as well as a medium for “fair exchange.”

257 Next we discuss the trust relationships that our system assumes. Firstly,
258 a data consumer trusts data providers to answer honestly. This is the hardest
259 trust assumption of our system, since due to the privacy-preserving properties
260 of our solution nobody can tell which providers responded honestly. This level
261 of trust can be (partially) achieved by implementing the provider-specific func-
262 tionality in a Trusted-Platform Module (TPM). Further discussion of such a
263 solution is out of the scope of this paper. Secondly, data consumers trust sys-
264 tem operators to perform correct filtering. Similarly, it should be impossible
265 for data consumers to verify that filtering was correctly applied, otherwise that
266 would mean that they have access to sensitive information. Nevertheless, since
267 the filtering rules and the filtering results are public, data providers can “fill a
268 complaint” or “rate a system operator negatively” in case they are incorrectly
269 excluded. Although we do not implement such mechanisms, in Section 5.3 we
270 present dispute resolution tools that are enabled by the use of the Ethereum
271 smart contract and can be used as building blocks for such systems. Thirdly,
272 data consumers and data providers trust system operators to distribute the
273 compensation paid by data consumers to data providers fairly. Again the tools

274 presented in section 5.3 can be used for building the appropriate reputation
275 systems.

276 In the next section we describe in the detail the design of our solution.
277 In order to better illustrate the concepts and their corresponding intuition we
278 consider the use case of a smart-grid measurements sharing system. In this use
279 case, the data providers are the deployed smart meters. A data consumer is a 3rd
280 party interested in extracting statistics about the grid, e.g., the average energy
281 consumption in a particular city. The system operator is the smart meters'
282 operator which holds information that is necessary for filtering the provided
283 responses. We assume that the system operator cannot provide (or is not allowed
284 to provide) the requested statistics.

285 4. Design and implementation

286 We assume that the system operator and the data providers are config-
287 ured with a pre-shared secret key (psk). Furthermore, they are configured with
288 a hash function $H(data)$, a keyed-hash message authentication code function
289 $HMAC(key, data)$, and a symmetric encryption algorithm $E(key, data)$. Each
290 data provider has an unlimited pool of Ethereum addresses and every time it
291 responds to a query it uses a different one: the system operator always knows
292 the current Ethereum address of each provider.

293 Data consumers should express the requested data in the form of a query with
294 multiple choices. For example, suppose a data consumer is interested in learning
295 the energy consumption per day of each smart meter, then a appropriate query
296 could be “What is your energy consumption per date? a) 1kW b) 2kW c) 3kW
297 ...”. No matter the number of responses, there is always the chance that a user
298 will respond randomly with probability 50%. Moreover, the number of responses
299 per choice that will be estimated by the data consumer is not affected by the
300 number of choices offered to data providers, hence the number of choices and
301 their granularity are determined based on the needs of each use case. Finally,
302 we record in the blockchain a cryptographic hash of a “configuration” file that

303 includes the query and the available choices, keeping this way the blockchain-
304 related cost fixed and independent of the number of choices. Additionally, a
305 data consumer can specify filtering rules for the specific query, e.g., “smart
306 meters must be located in Athens, Greece.” The protocols implemented by our
307 architecture are discussed next.

308 4.1. Smart contract creation

309 With this protocol, a system operator and a data consumer agree on a “sur-
310 vey configuration” that includes (i) the query $Q = \{c_1, c_2, \dots, c_n\}$ that data
311 providers should respond, (ii) the filtering rules, and (iii) the number of re-
312 sponses N_R that should be collected. The location of this configuration (e.g.,
313 a URL) as well as its hash are stored in a smart contract. Additionally, the
314 system operator and the data consumer agree on two nonces $n1$ and $n2$. All
315 data providers and the system operator (but not the data consumers) can de-
316 rive an encryption key $sk = H(s1||s2)$, where $s1 = HMAC(psk, n1)$, $s2 =$
317 $HMAC(psk, n2)$, and $||$ denotes concatenation. Finally, $s1$ is securely transmit-
318 ted to the data consumer, and $n1$, $n2$, $H(s2)$, and a “service fee” are recorded
319 in the smart contract.

320 4.2. Response commit

321 Data providers can retrieve the survey configuration from the specified loca-
322 tion and verify its integrity. Any data provider wishing to respond to a query,
323 prepares a “noisy” response R using local differential privacy. The response is
324 constructed using the “basic one-time RAPPOR” algorithm (see section 2.2).
325 In particular, R is a bit vector of size equal to the query Q . The value of each
326 element $r_i \in R$ is equal to the output of the randomized response game (also
327 discussed in Section 2.2) for the choice $c_i \in Q$, i.e., the value of the first ele-
328 ment of R is the output of the randomized response game for the first choice in
329 Q , and so forth. Then, the data provider derives the encryption key sk using
330 the nonces $n1$ and $n2$ included in the smart contact and generates a ciphertext
331 $C_R = E(sk, R)$, i.e., it encrypts the bit vector R with sk . Finally, it records

332 the hash of the generated ciphertext $H(C_R)$ in the smart contract. All hashes
333 $H(C_R)$ are stored in the smart contract in a sorted map $M_{address \rightarrow H(C_R)}$ that
334 maps the Ethereum address of the data provider to the corresponding $H(C_R)$.

335 4.3. Response filtering

336 A system operator maintains a filter F . F is a dynamic bit vector whose
337 size is increased by one every time an $H(C_R)$ is recorded in the map M of the
338 smart contract. The value of a bit $f_i \in F$ is set to 1 if the i^{th} data provider
339 that recorded a response in the smart contract meets the agreed filtering criteria
340 otherwise it is set to 0. The implementation of the process for deciding about
341 whether a response meets the filtering criteria is application specific. When the
342 number of ones in F is equal to N_R (i.e., the number of the responses that
343 can be accepted equals to the number of the required responses) the system
344 operator generates the hash $H(F)$, records it in the smart contract, and makes
345 F available to the data providers and to the data consumer (e.g., it publishes
346 in a pre-agreed URL). At this point no new responses are accepted.

347 4.4. Response submission

348 All data providers that have committed a response (by submitting $H(C_R)$
349 to the smart contract) retrieve the filter F (from the system operator), as well
350 as the map M (from the smart contract). Based on F a data provider can
351 learn if it has been accepted to submit its response by executing the following
352 process. It locates the position i in M that corresponds to its Ethereum address,
353 and it examine the value of $f_i \in F$ (i.e., the i^{th} element of F). If $f_i = 1$ then
354 it is accepted and it sends i and the corresponding ciphertext C_R to the data
355 consumer (it is reminded that $C_R = E(sk, R)$, i.e., C_R is the encryption of the
356 noisy response using sk).

357 4.5. Fair exchange

358 Upon receiving i and C_R from a provider, the data consumer verifies (a)
359 the integrity of C_R using the hash $H(C_R)$ stored in the i^{th} position of the

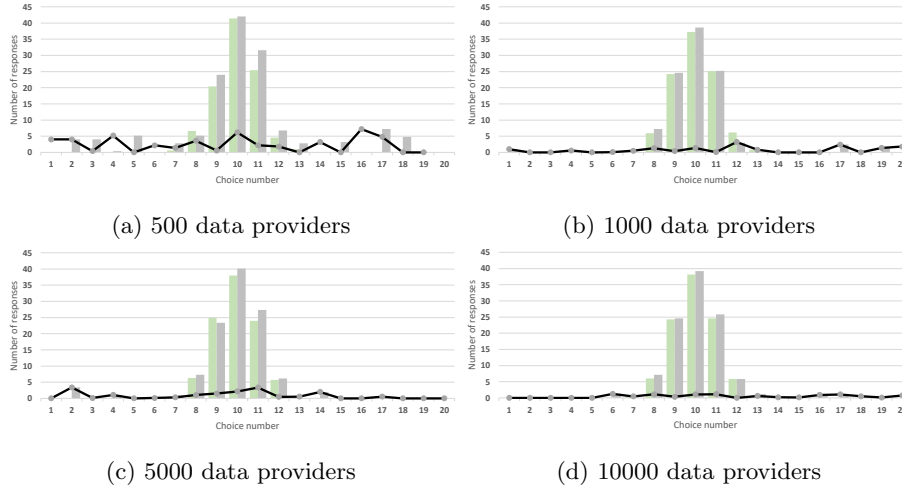


Figure 3: The impact of local differential privacy. Green axes show the real number of responses. Grey axes the estimated responses. The black line shows their difference: the closer the black line is to x axis the better.

360 map M (retrieved from the smart contract), and (b) that the data provider
 361 has been accepted to submit C_R by examining if the value of the i^{th} element
 362 of F (retrieved from the system operator) is 1. When it receives the agreed
 363 number of responses it deposits to the smart contract the appropriate amount
 364 of currency. Then the system operator records $s2$ to the smart contract; if the
 365 hash of the revealed $s2$ matches the hash $H(s2)$ stored in the smart contract,
 366 the contract transfers the deposit to the system operator. The system operator
 367 is then trusted to compensate the proper data providers accordingly. The data
 368 consumer can now reconstruct the data encryption key using $s1$ (received with
 369 the smart contract creation protocol), and $s2$, simply by calculating $H(s1||s2)$.

370 4.6. Results extraction

371 The data consumer constructs sk and uses it to decrypt the received C_R 's.
 372 Each decrypted message is the noisy response of a data provider. Then, using
 373 the formula of Section 2.2 it calculates the probability of each choice $c_i \in Q$.

374 5. Evaluation

375 5.1. Local differential privacy efficiency

376 The accuracy of RAPPOR (and of local differential privacy in general) de-
377 pends on the number of the submitted responses. In order to evaluate the
378 accuracy of the extracted statistics, we perform the following experiment.¹ We
379 create a number of data provider instances (the number of instances is a variable
380 in our experiments) and we ask them to respond to a query that includes 20
381 possible choices. Each provider instance selects its response using a normal dis-
382 tribution with mean value 10 and deviation 2. Then each data provider submits
383 its response. We then plotted (Figure 3) the real distribution of the responses
384 (green bar), the extracted distribution based on the noisy responses (black bar),
385 and their difference (black line): the closer to the horizontal axis the black line
386 is, the more accurate the results are. The following diagrams concern experi-
387 ments with 500, 1000, 5000, and 10000 data providers. It should be noted that
388 the accuracy of the extracted results is not affected by the number of choices,
389 neither by the distribution of the real responses.

390 5.2. Blockchain-based overhead

391 We implemented² a smart contract that provides the functionality described
392 in Section 4 and we deployed and tested it in the Rinkeby Ethereum test net-
393 work.³ The following table shows the cost measured in units of “gas” for de-
394 ploying our smart contract, as well as for invoking its functions.

395 The fiat cost included in this table is based on the prices retrieved from
396 <https://ethgasstation.info> on 1 Jun. 2021. The cost of each operation is not
397 affected by the number of data providers neither by the number of available
398 choices.

¹Our implementation, as well as instructions for replicated the experiments can be found
at <https://github.com/mmlab-aueb/rappor>

²The implemented system can be found in <https://github.com/mmlab-aueb/Privacy-and-Data-Sovereignty>

³<https://www.rinkeby.io/>

Table 1: Smart contract costs

Operation	Cost measured in gas	Cost measured in fiat
Contract Deployment	660809	\$32,5
Record $H(C_R)$	74537	\$3,66
Record $H(F)$	63309	\$3,11
Make deposit	23642	\$1,16
Reveal s_2 and transfer deposit	36269	\$1.78

399 5.3. Dispute resolution

400 A key role of the smart contract in our system is to act as an immutable,
 401 append-only log where the outcomes of a number of hash functions are recorded.
 402 These records can be latter used with a dispute resolution mechanism.

403 Each data provider records in the smart contract $H(C_R)$. Then data providers
 404 send C_R to data consumers. In case C_R is unreadable, e.g., it cannot be de-
 405 crypted with sk , the data consumer can use $H(C_R)$ in a dispute resolution
 406 process: since $H(C_R)$ is recorded in the blockchain, a data provider cannot
 407 deny it sent C_R .

408 Similarly, system operators record $H(F)$. A data consumer not included in
 409 the filter F can file a complain. Since $H(F)$ is recorded in the blockchain a
 410 system operator cannot deny that it used F . $H(F)$ can also be used by dispute
 411 resolution system that handles cases where a data provider is included in F but
 412 it did not receive the appropriate compensation.

413 5.4. Privacy properties of RAPPOR

414 We now discuss the privacy properties of RAPPOR. We consider the case of
 415 a question that includes n choices, i.e., $Q = \{c_1, c_2, \dots, c_n\}$ and the real responses
 416 follow a uniform distribution, hence for every provider $P_{(x=C)} = 1/n$, where C is
 417 the choice of the provider. The probability that an attacker with no additional
 418 information guesses the choice C of a provider is $P_{guess} = 1/n$. We now consider
 419 an attacker that can access a randomized response R of a provider, generated

420 using RAPPOR. The attacker tries to guess C by selecting at random an element
 421 of $r_i \in R$ whose value is 1. The probability that $r_i = 1$ is given by calculating
 422 the following formula.

$$P_{(r_i=1)} = \frac{1}{n} \times 0.5 + 0.25 \quad (1)$$

423 This happens because an element is 1 either if it corresponds to the actual
 424 choice C of the provider and the output of the randomized game is *true* (which
 425 has probability $\frac{1}{n} \times 0.5$) or the output of the randomized game is “say 1” (which
 426 occurs with probability 0.25). Therefore the numbers of 1s in R is $\approx P_{(r_i=1)} \times n$.
 427 Furthermore, the probability that an element of this set is C is 75%. We define
 428 the probability that this attacker guess correctly C as P_{RAPPOR} and we define
 429 the advantage of this attacker as $Adv = P_{RAPPOR}/P_{guess}$. Figure 4 shows these
 430 probabilities and the advantage of the attacker which approximates 3. This
 431 means that even though the P_{RAPPOR} decreases as the number of responses
 432 increases, it is ≈ 3 times bigger than P_{guess} . An intuitive explanation for this is
 433 that as the number of choices increases, $P_{(r_i=1)}$ approximates 0.25, hence, the
 434 number of 1s in R approximates $0.25 \times n$ therefore, $P_{RAPPOR} \approx \frac{1}{0.25 \times n} \times 0.75 =$
 435 $\frac{3}{n} = 3 \times P_{guess}$

436 An interesting trade-off that we can consider is to assume a “non-fair” coin
 437 for the first round of the randomized response game. For example a coin that
 438 decides that a user will say the truth 20% of the times (as opposed to 50%)
 439 decreases the advantage of the attacker to 1.5 (since, as the number of choices
 440 increases, the number of 1s in R approximates $(0.8 \times 0.5) \times n$ and the probability
 441 that an element of this set is C is 60%). Figure 5 shows the difference between
 442 real distributions of the responses and the distribution based on the noisy re-
 443 sponses for 10000 data providers and 20 choices, selected using the distribution
 444 presented in section 5.1, using a fair coin (black line) and a coin that decides
 445 that the user will say the truth 20% of the times. As it can be seen, by using
 446 the second coin we reduce the advantage of the attacker, but at the cost of lower
 447 accuracy.

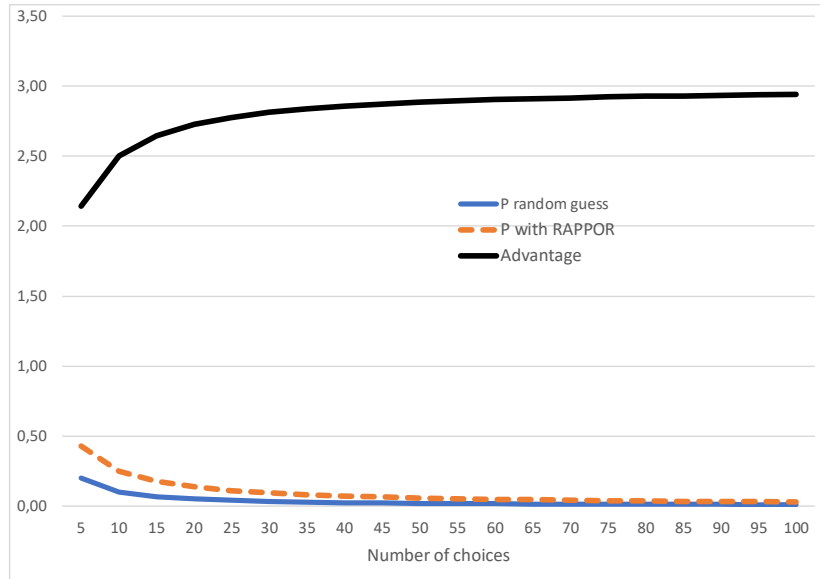


Figure 4: Probability of a successful random guess, and of a guess based on a RAPPOR response. The “Advantage” line shows the advantage of an attacker that has access to a RAPPOR response.

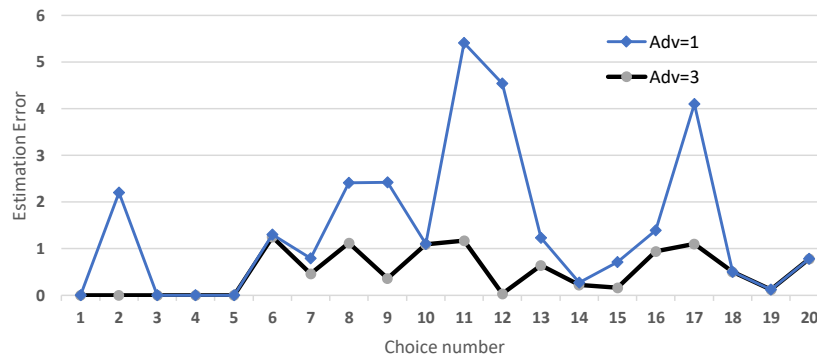


Figure 5: Error in estimation when different levels of attacker’s advantage are considered.

448 5.5. Privacy-accuracy trade-off

449 The threat model considered in differential privacy mechanisms is that of
450 an attacker that can monitor how the computed statistic “evolves” with the
451 addition (or removal) of samples. Solutions that provide accurate statistics
452 fail to provide protection against this type of attackers. In order to illustrate
453 this attack, we perform the following experiment: we consider providers that
454 select uniformly and at random a choice from the set $0, 1, \dots, 19$ and securely
455 commit their response to a trusted service; every time a provider commits a
456 response, the service outputs the average of all committed responses. After
457 1000 providers have committed their responses, an attacker starts monitoring
458 the extracted average. Then, a new provider commits a response, the new
459 average is calculated, and the attacker tries to guess based on the new average
460 what was the choice of the provider. We consider two cases: (a) the responses
461 are committed without noise addition (e.g., a mechanism based on homomorphic
462 encryption is used), hence the extracted results are accurate, and (b) noise is
463 added to the responses using RAPPOR. We perform the same experiment with
464 all possible choices the 1001st provider can make. Figure 6 shows the choice
465 estimated by the attacker. The horizontal axis of the diagram included in this
466 figure corresponds to the correct choice. As it can be seen, when no noise is
467 added the attacker can successfully guess the added choice. However, when
468 RAPPOR is used this is not possible.

469 6. Discussion

470 Our system has the following properties:

471 **It is fast.** Data providers have only to record a single hash in the smart
472 contract and then send a small message (in order of hundreds of bytes).

473 **It is lightweight.** The most computationally intensive operation that a
474 data provider has to perform is the calculation of a hash, the calculation of a
475 digital signature, and a symmetric key encryption. Similarly, a data consumer
476 has only to perform hash calculations and multiplications, and a symmetric key

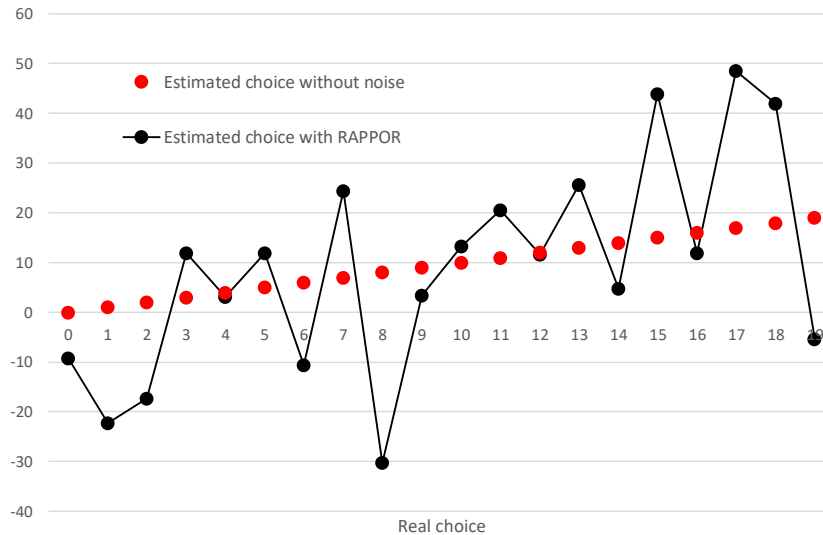


Figure 6: Gussed value when no noise is added to the response (red dots) and with RAPPOR (black dots). The horizontal axis corresponds to the correct choice.

477 decryption. All these computations can be performed fast even by constrained
 478 devices.

479 Many of the design choices in our system are driven by the need for support-
 480 ing many data providers. For example, having the system operator distributing
 481 the service fee to the data providers, instead of the smart contract itself, may
 482 appear strange. However, implementing digital currency transfer to hundreds
 483 of users in a smart contract has prohibitive cost (and if the number of users is
 484 very big such an operation may not even be allowed). An alternative approach
 485 can be the implementation of a “withdrawal” function in the smart contract
 486 that eligible data providers can invoke, after the release of s_2 , and retrieve their
 487 compensation by themselves. Similarly, one may argue that the smart contract
 488 should filter the map M based on the filter provided by the system operator,
 489 nevertheless since M contains hundreds of records this is a very costly operation
 490 to implement in a smart contract.

491 A parameter of our system—as well as of any RAPPOR-based system—is the
 492 number of available choices per question. We discuss in Section 5.4 that, as the

493 number of choices increases, the advantage, Adv , of an attacker that has access
494 to a RAPPOR response remains stable. Of course, the bigger the number of
495 choices the better the privacy protection for providers since the probability of
496 a successful guess is $Adv \times P_{guess}$, where P_{guess} is the probability of randomly
497 guessing the choice of a provider; as the number of choices increases P_{guess}
498 decreases. Moreover, as the number of choices increases the accuracy of the
499 extracted results also increases. Nevertheless, and as we show in Section 5.1, the
500 use of RAPPOR introduces some error in the extracted results which depends on
501 the number of providers, e.g., it can be up to 5% with 500 providers and close to
502 1% with 10,000 providers. Therefore, the number of choices should be selected
503 such that the distribution of the extracted results makes sense even with the
504 error rate of RAPPOR. For example, in a scenario with 500 providers, offering
505 many choices would result in a distribution where the most popular choices are
506 selected by $\approx 5\%$ of the providers, will make the real results indistinguishable
507 from the error rate.

508 Our system can be modified to support even higher privacy or greater de-
509 centralization. The former case is related to the fact that a data consumer can
510 perform some form of data provider tracking based on the provider’s network
511 address (recall that the Ethereum address of each provider changes whenever
512 they commit a response). Therefore, an alternative could be that the data
513 providers send their encrypted responses to the system operator, and then the
514 system operator to forward all of them to the data consumer. On the other
515 side, it can be argued that some type of filtering can be performed by the data
516 consumer (or even the smart contract) if data consumers can prove some of their
517 properties (e.g., using *Verifiable Credentials* [26] that include their location). In
518 that case the role of the system operator would be reduced and limited, e.g., to
519 issuing the VC for verifying the data consumer’s location or any other property
520 used for filtering.

521 Local differential privacy systems—such as ours—are susceptible to “longitu-
522 dinal” attackers, i.e., attackers that have access to multiple reports over time
523 from the same user. RAPPOR protects users from this attack by introducing a

524 “memoization” step [9]. With this step, when a user is offered the same choice,
525 the user responds with the same, permanently stored response. Nevertheless,
526 defining whether two choices are the same may require human intervention, e.g.,
527 the choices “consumption at night is between 1kW and 5kW”, and “consump-
528 tion from 6pm to 4am is between 1kW and 5kW” can be considered the same.
529 We postulate that since the service provider knows all the questions that have
530 been asked it is its role to inform data consumers if a choice is the same to one
531 previously offered.

532 Another aspect that affects the deployability of our system is the cost of
533 transacting with the public blockchain. As it can be seen from Table 1 currently
534 the most expensive operation (apart from deploying a smart contract) has cost
535 \$3,66. These costs may limit the applicability of our approach for certain use
536 cases. Similarly, these costs have high fluctuation, which is another problem in
537 itself, since they depend on the price of ETH, i.e., the Ethereum digital currency.

538 **7. Conclusions and future work**

539 In this paper we presented a privacy-preserving marketplace for sensitive
540 data. Our solution allows data consumers to “purchase” noisy data that can
541 be used for extracting meaningful statistics. The privacy of data providers is
542 protected against all entities of our system using local differential privacy. Our
543 solution allows intermediate service providers that can filter out data providers
544 without having access to the provided data neither to the extracted statistics.
545 Our solution uses a blockchain-based approach for providing an immutable log,
546 as well as for implementing fair exchange. Although our system may involve
547 many stakeholders, the blockchain-based overhead is small and independent of
548 the number of data providers.

549 The use of local differential privacy means that nobody can tell if a data
550 provider responded to a query honestly. Therefore, increased privacy comes at
551 the cost of exposure to malicious data providers. An interesting approach for
552 mitigating this issue is the execution of the data provider specific functionality

553 in a Trusted Platform Module (TPM): the use of TPMs will be considered in
554 our future work.

555 Ethereum smart contracts enable fair exchange, as well as immutable logs.
556 Although our design keeps interactions with the blockchain to a minimum, and it
557 has a fixed blockchain-based cost, no matter the number of the stakeholders, the
558 overhead, the complexity, and the monetary cost introduced by this technology
559 cannot be ignored. To this end, we are studying alternatives based on private,
560 consortium-based blockchain systems.

561 **Acknowledgment**

562 This research was supported by the EU funded Horizon 2020 project SOFIE
563 (Secure Open Federation for Internet Everywhere), under grant agreement No.
564 779984.

565 **References**

- 566 [1] D. Susser, B. Roessler, H. Nissenbaum, Online manipulation: Hidden influ-
567 ences in a digital world, *Georgetown Law Technology Review*, Forthcoming.
- 568 [2] Z. J. T. Barbaro M., A face is exposed for aol searcher no. 4417749 (2006).
569 URL <https://www.nytimes.com/2006/08/09/technology/09aol.html>
- 570 [3] A. Narayanan, V. Shmatikov, Robust de-anonymization of large sparse
571 datasets, in: *2008 IEEE Symposium on Security and Privacy (SP 2008)*,
572 2008, pp. 111–125.
- 573 [4] P. Ohm, Broken promises of privacy: Responding to the surprising failure
574 of anonymization, *UCLA L. Rev.* 57 (2009) 1701.
- 575 [5] C. Dwork, F. McSherry, K. Nissim, A. Smith, Calibrating noise to sensi-
576 tivity in private data analysis, in: S. Halevi, T. Rabin (Eds.), *Theory of*
577 *Cryptography*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2006, pp.
578 265–284.

- 579 [6] E. C. D. Energy, Benchmarking smart metering deployment in the eu 28,
580 revised final report (2019).
581 URL [https://www.vert.lt/SiteAssets/teises-aktai/EU28%20Smart%
582 20Metering%20Benchmark%20Revised%20Final%20Report.pdf](https://www.vert.lt/SiteAssets/teises-aktai/EU28%20Smart%20Metering%20Benchmark%20Revised%20Final%20Report.pdf)
- 583 [7] S. Dziembowski, L. Eckey, S. Faust, Fairswap: How to fairly exchange
584 digital goods, CCS '18, Association for Computing Machinery, New York,
585 NY, USA, 2018, p. 967–984.
- 586 [8] S. L. Warner, Randomized response: A survey technique for eliminat-
587 ing evasive answer bias, *Journal of the American Statistical Association*
588 60 (309) (1965) 63–69.
- 589 [9] U. Erlingsson, V. Pihur, A. Korolova, Rappor: Randomized aggregatable
590 privacy-preserving ordinal response, in: *Proceedings of the 2014 ACM*
591 *SIGSAC Conference on Computer and Communications Security, CCS '14,*
592 *ACM, New York, NY, USA, 2014, pp. 1054–1067.*
- 593 [10] F. McSherry, I. Mironov, Differentially private recommender systems:
594 Building privacy into the netflix prize contenders, in: *Proceedings of the*
595 *15th ACM SIGKDD International Conference on Knowledge Discovery and*
596 *Data Mining, KDD '09, Association for Computing Machinery, New York,*
597 *NY, USA, 2009, p. 627–636.*
- 598 [11] A. Friedman, A. Schuster, Data mining with differential privacy, in: *Pro-*
599 *ceedings of the 16th ACM SIGKDD International Conference on Knowl-*
600 *edge Discovery and Data Mining, KDD '10, Association for Computing*
601 *Machinery, New York, NY, USA, p. 493–502.*
- 602 [12] J. Hamm, A. C. Champion, G. Chen, M. Belkin, D. Xuan, Crowd-ml: A
603 privacy-preserving learning framework for a crowd of smart devices, in:
604 *2015 IEEE 35th International Conference on Distributed Computing Sys-*
605 *tems, 2015, pp. 11–20.*

- 606 [13] A. Mani, M. Sherr, *Histor ϵ : Differentially private and robust statistics*
607 *collection for tor.*, in: NDSS, 2017.
- 608 [14] F. Kargl, A. Friedman, R. Boreli, *Differential privacy in intelligent trans-*
609 *portation systems*, in: Proceedings of the Sixth ACM Conference on Secu-
610 *rity and Privacy in Wireless and Mobile Networks, WiSec '13*, Association
611 *for Computing Machinery, New York, NY, USA, 2013*, p. 107–112.
- 612 [15] A. Friedman, I. Sharfman, D. Keren, A. Schuster, *Privacy-preserving dis-*
613 *tributed stream monitoring.*, in: NDSS, 2014, pp. 1–12.
- 614 [16] T. Dimitriou, A. Mohammed, *Fair and privacy-respecting bitcoin payments*
615 *for smart grid data*, IEEE Internet of Things Journal 7 (10) (2020) 10401–
616 10417. doi:10.1109/JIOT.2020.2990666.
- 617 [17] H. Duan, Y. Zheng, Y. Du, A. Zhou, C. Wang, M. H. Au, *Aggregating*
618 *crowd wisdom via blockchain: A private, correct, and robust realization*,
619 *in: 2019 IEEE International Conference on Pervasive Computing and Com-*
620 *munications (PerCom, 2019)*, pp. 1–10.
- 621 [18] M. Burkhart, M. Strasser, D. Many, X. Dimitropoulos, *Sepia: Privacy-*
622 *preserving aggregation of multi-domain network events and statistics*, in:
623 *Proceedings of the 19th USENIX Conference on Security, USENIX Secu-*
624 *rity'10*, USENIX Association, USA, 2010, p. 15.
- 625 [19] G. Ács, C. Castelluccia, *I have a dream! (differentially private smart meter-*
626 *ing)*, in: T. Filler, T. Pevný, S. Craver, A. Ker (Eds.), *Information Hiding*,
627 *Springer Berlin Heidelberg, Berlin, Heidelberg, 2011*, pp. 118–132.
- 628 [20] M. Hale, P. Barooah, K. Parker, K. Yazdani, *Differentially private smart*
629 *metering: Implementation, analytics, and billing*, in: Proceedings of the 1st
630 *ACM International Workshop on Urban Building Energy Sensing, Controls,*
631 *Big Data Analysis, and Visualization, UrbSys'19*, Association for Comput-
632 *ing Machinery, New York, NY, USA, 2019*, p. 33–42.

- 633 [21] C. Liu, S. Chen, S. Zhou, J. Guan, Y. Ma, A general framework for privacy-
634 preserving of data publication based on randomized response techniques,
635 Information Systems 96 (2021) 101648.
- 636 [22] K. Gai, Y. Wu, L. Zhu, Z. Zhang, M. Qiu, Differential privacy based
637 blockchain for industrial internet-of-things, IEEE Transactions on In-
638 dustrial Informatics 16 (6) (2020) 4156–4165. doi:10.1109/TII.2019.
639 2948094.
- 640 [23] F. Fioretto, T. W. K. Mak, P. Van Hentenryck, Differential privacy for
641 power grid obfuscation, IEEE Transactions on Smart Grid 11 (2) (2020)
642 1356–1366.
- 643 [24] J. Zhao, T. Jung, Y. Wang, X. Li, Achieving differential privacy of data
644 disclosure in the smart grid, in: IEEE INFOCOM 2014 - IEEE Conference
645 on Computer Communications, 2014, pp. 504–512.
- 646 [25] M. Ul Hassan, M. H. Rehmani, J. Chen, Differential privacy in blockchain
647 technology: A futuristic approach, Journal of Parallel and Distributed
648 Computing 145 (2020) 50–74.
- 649 [26] Manu Sporny et al., Verifiable credentials data model 1.0 (2019).
650 URL <https://www.w3.org/TR/verifiable-claims-data-model/>