# Privacy-Preserving Release of Energy Data

July 2019

# Summary

As electricity utility proceed with the replacement of conventional electricity meters by deploying smart electricity meters in millions of home around the world, an immense amount of fine-grained electricity consumption data are being automatically and frequently collected. Utilising the big data available in such energy data collection, a variety of data analytics algorithms and applications have been developed. However, growing privacy concerns and confidentiality issues preclude the sharing or releasing of energy data that contain or might enable someone to deduce personal information (such as occupancy rate and life style) [Molina-Markham et al., 2010, Wood and Newborough, 2003, Wang and Zheng, 2012, McDaniel and McLaughlin, 2009, Greveler et al., 2012]. This reduces the ability to use such data for data analytics and applications. This report studies the research problem in two directions. The first one is privacy-preserving release of percentile statistics for energy data and the second one is privacy preserving clustering of energy consumption data.

In the first part, we develop privacy-preserving policies for releasing percentile statistics, e.g., median, of time-series data. Although percentile statistics are aggregating information, they are not privacy-preserving. For instance, the median of a vector with an odd number of elements is equal to the value of the middle entry of the vector after being sorted. The privacy of that middle entry is not respected. The privacy concerns for the middle entry, whichever it is, are not alleviated even if the number of data entries grow. We base our analysis on the research areas of differential privacy, local differential privacy, stochastic differential privacy, and information theory. We prove that for percentile statistics, local differential privacy and stochastic differential privacy outperform differential privacy in terms of the quality of reported output (i.e., "closeness" of the reports to the desired percentile statistics). Local differential privacy is also preferred as it provides privacy guarantees in the presence of an untrusted actor, such as third party data aggregators. In addition, the privacy guarantees of the stochastic differential privacy is comparatively weaker as it enforces the privacy requirements on almost all datasets rather than all datasets. Finally, we turn our attention to information-theoretic privacy and develop optimal privacy-preserving policies by maximizing the entropy of the additive noise subject to constraints on various measures of quality corresponding to expectation of the absolute value, variance, and essential entry-wise extrema of the additive noise. We use the publicly available dataset of half-hour energy measurements from 300 homes in Australia with rooftop solar panels to demonstrate the results.

In the second part, we investigate the privacy risks of clustering household consumption data as well as the use of local differential privacy for extracting privacy-preserving policies. We first remove each house within the dataset to observe the new clusters. Noting that the clusters change massively, up to 2% of the households change clusters, there might be a risk of privacy infringement from releasing the clusters without privacy treatments. To alleviate the risk, we propose the use of local differential privacy. The differential privacy noise can potentially change the clusters massively for reasonable choices of differential privacy parameters. This further points to the sensitive nature of the clusters and their possible privacy risks.

The results indicate that the local differential privacy outperforms other methods for privacy-preserving releasing of percentile statistics in terms of the privacy-utility trade-off. Differential privacy also removes the need for trusted third-party for computations. Moreover, the investigation indicate that releasing clusters of household consumption data can compromise privacy and the use of local differential privacy improves privacy-preservation at the cost of significant utility loss. Hence, we require further research in the topic of privacy-preserving clustering of energy consumption data.

# Contents

# 1 Introduction

Smart meters are gradually replacing traditional electricity meters enabling remote collection of high-frequency energy data. While the primary purpose of the data is to enable time of day billing of electricity usage, the collected data provides the opportunity for analysis to pave the way for governments, energy providers, and grid operators to make informed decisions resulting in better quality or service, improved reliability, and lower consumer prices; see, e.g., [Pappu et al., 2017, Perez et al., 2017, Wang et al., 2018]. Privacy concerns[1], however, prohibit the release of raw smart meter data as it reveals potentially sensitive information about the participants, such as occupancy rate, domestic appliance usage, dietary habits, and even multimedia consumption [Molina-Markham et al., 2010, Wood and Newborough, 2003, Wang and Zheng, 2012, McDaniel and McLaughlin, 2009, Greveler et al., 2012]. In this report, we address two problems related to privacy-preserving release of smart meter energy data. The former is privacy-preserving release of percentile statistics of the energy data, and the latter is privacy-preserving clustering of household consumption data.

Most often, governments and private entities are interested in aggregate statistics of energy consumption to inform decisions on load balancing and future investment. A class of highly sought-after statistics are percentile statistics, such as median, as they provide desirable properties, such as robustness to outliers or malicious data injections [Tukey et al., 1983, Farokhi et al., 2015b, Tukey, 1960]. Although providing aggregate information, percentile statistics are not privacy preserving. This is because the percentile statistics of a vector (that is a series of numbers) correspond to one or two entries of the vector. For instance, the median of a vector is equal to the value of the middle entry of the sorted vector if the number of the entries is odd or the average of the two middle entries of the sorted vector if the number of entries is even. Therefore, as the number of the participants grow (i.e., the dataset expands), the privacy concerns do not get alleviated as one (or two) entries are exposed. (This is contrary to other means of aggregation, such as arithmetic mean). Previously, other studies have developed privacy-preserving policies for energy data at an aggregate level using differential privacy [Eibl and Engel, 2017, Sandberg et al., 2015]. However, those studies have not considered percentile statistics. They also do not investigate the performance of other relevant notions of privacy, such as local differential privacy and stochastic differential privacy; see below.

One of the most prominent privacy definitions for providing privacy guarantees is differential privacy [Dwork, 2008, Dwork and Roth, 2014, Dwork et al., 2010]. It guarantees that, given the published output, any information that could be discovered about an individual with their data in the underlying dataset could also, with high probability, be discovered without their data in the dataset. However, regardless of the amount of auxiliary information the adversary has about an individual, it is unable to identify the individual's presence in the dataset with high probability

---

[1]Such concerns have proved to even hinder the roll out of smart meters in some countries [Cuijpers and Koops, 2013].

(determined by a privacy parameter). In other words, a differentially-private algorithm perturbs the output by adding noise in a way that protects the privacy of every individual in the original dataset. When perturbing the percentile statistics with an additive noise following differential privacy definition, it guarantees that the distribution of the reported percentile statistics does not significantly change when an individual data entry in the underlying dataset changes (for example, the energy data time series of an individual changes). Perturbing output to provide privacy guarantees obviously incurs utility loss, i.e. the deviations of the reported outputs from the true values of percentile statistics (essentially a measure of the magnitude of the additive noise, such as its variance). However, as discussed in detail in Section 1.1, the privacy-utility trade-off of the perturbed output with differentially-private additive noise is often unappealing for percentile statistics output in contrast to the output of arithmetic mean or counts.

An alternative approach to ensure privacy in the reported percentile statistics output is local differential privacy [Dewri, 2013, Duchi et al., 2013, Kairouz et al., 2014]. Local differential privacy ensures that the data is kept private from the aggregator by adding noise to the individual data entries before the aggregation process rather than the aggregated output itself. This is a preferred choice when dealing with untrusted aggregators, e.g., third party service providers with financial interests in the energy data. For percentile statistics, we prove that local differential privacy provides a better utility-privacy trade-off. That is, fixing the privacy requirement, the utility of local differential privacy is higher. This makes the local differential privacy the method of choice even if the aggregator is trusted due to its superior utility. A more relaxed privacy-preserving approach is to use stochastic differential privacy [Machanavajjhala et al., 2008, Rubinstein and Aldà, 2017, Hall et al., 2012] which requires the definition to be held for almost all datasets, but not all. This is in contrast to the all datasets definition with the differential privacy definition. With this approach the utility loss of perturbed output reduces at the cost of weaker privacy guarantees due to the relaxation of the definition.

A very different approach to the design of privacy-preserving policies is to use information-theoretic tools for capturing information leakage and privacy infringements [Wang et al., 2016, Tanaka et al., 2018, 2017, Kalantari et al., 2017, du Pin Calmon and Fawaz, 2012]. Examples of various metrics for capturing information leakage are the least mean square error [Farokhi et al., 2015a], mutual information or relative entropy [Tan et al., 2013], directed information [Tanaka et al., 2017], hypothesis testing error rate [Li et al., 2017], and Fisher information [Farokhi and Sandberg, 2017]. Information-theoretic measures of privacy capture the behaviour of the percentile statistics and privacy-preserving additive noises in a statistical average sense and are thus immune to overreacting to the aforementioned property of percentile statistics that can only depend on one data entry (in extreme, yet unlikely cases).

This report presents our methods, theoretical proofs, and findings of privacy-preserving release of percentile statistics using different approaches (as briefly described above), especially when applied to smart meter energy data. Moreover, clustering household energy consumption data with the aim to provide data analysts and researchers useful insights about the data, might have the risk of revealing private or sensitive information about individuals in

the data, and at least has a theoretical risk. A privacy-preserving clustering algorithm requires that the reported clusters do not change noticeably when an individual's data is removed from the underlying dataset. In this study, we evaluate the privacy risk of releasing clusters of data and the privacy-utility trade-off of privacy-preserving clustering using local differential privacy.

**Contributions:** Our contributions are two-fold. First, we contribute to the problem of privacy-preserving release of percentile statistics (median) of smart meter energy data in Chapter 2 and then we study the problem of privacy-preserving clustering of household consumption data in Chapter 3. We use different notions of privacy to address these problems and we evaluate and compare the privacy-utility trade-off using publicly available real datasets.

## 1.1     Privacy-Preserving Release of Percentile Statistics of Energy Data

In the following chapter, we tackle the problem of privacy-preserving release of percentile statistics from multiple angles:

- *Differential Privacy*: We first develop differentially private policies for the release of percentile statistics. Universal mechanism, such as Laplace, Gaussian, and exponential mechanisms [Dwork and Roth, 2014], have proved to be easy to implement and popular for ensuring differential privacy. The Laplace (Gaussian) mechanism relies on systematically corrupting the output using an additive Laplace (Gaussian) noise whose scale (variance) is proportional to the sensitivity of the reporting function to variations of individual entries in the dataset. Noting that smart meters (or other energy data meters) report time series, we investigate various notions of adjacency for the datasets. The notion of adjacency is the basis of computing sensitivity of reporting function to variations of individual entries of the dataset. We consider the following notions of adjacency for differential privacy:

  - *Point-wise Adjacency*: This is a straightforward extension of adjacency from static datasets to time series. The point-wise adjacency states that two datasets are adjacent if they differ from each other in one time instance and for at most one time series. Essentially, in this case, we simplify the time series problem to multiple independent static snapshots of it (one snapshot for each time instance of the time series).

  - *Trajectory Adjacency*: Energy data time series are outcomes of physical systems pointing to their overall smoothness. This motivates us to trajectory adjacency, which states that two datasets are adjacent if they differ from each other for at most one time series and, for those time series, the values remain within a reasonable distance of each other. This notion of adjacency is highly relevant when the time series follows a discrete-time difference equation (i.e., capturing an automata) and the difference between the time series originates from the difference in the initial conditions of the difference equations, perturbations to the parameters of the

automata, or exogenous inputs to the automata.

The sensitivity of percentile statistics to variations of individual entries in the dataset is large and does not decrease with the size of the dataset (contrary to the arithmetic mean whose sensitivity is inversely proportional to the size of the dataset). This is because of the existence of scenarios in which the percentile statistics depend only on the data of one individual. For instance, consider the median of a vector with an odd number of entries such that half of the entries are equal to the minimum allowed value and the other half are equal to the maximum allowed value. Now, the median is equal to the last remaining entry and thus the sensitivity of the median to the changes of one entry is as large as the range of allowed values (i.e., the changes of the free entry).[2] Intuitively, in such scenarios, the median is not aggregative and infringes on the privacy of that individual (by directly reporting its private energy data) and, therefore, we need to utilize a fairly large additive noise to mask the private information of the-said individual. This makes the privacy-utility trade-off unappealing for percentile statistics with differentially-private additive noise. Here, utility is a measure of the deviations of the reported outputs from the queried percentile statistics (essentially a measure of the magnitude of the additive noise, such as its variance). Such pathological cases (i.e., the cases in which the percentile statistics can arbitrarily change by varying the data of one individual) are not likely to occur if the underlying dataset is probabilistically distributed (and its probability density function is atom-less[3]). Nonetheless, differential privacy does not assume a prior on the underlying dataset and must deal with such "pathological" cases. So we require an additive noise with a large magnitude to meet the differential privacy definition (which will undermine the utility).

- *Local Differential Privacy*: We also develop reporting policies for the release of percentile statistics based on local differential privacy. For local differential privacy, in addition to local reformulation of *point-wise adjacency* and *trajectory adjacency*, we consider the following notion of adjacency:

  - *Sparse Adjacency*: Time series associated with physical systems are often sparse in another domain, such as the frequency domain. This prompt us to investigate adjacency within a sparsifying domain. The sparse adjacency states that two datasets are adjacent if they differ from each other for at most one time series and, for those time series, the values within the sparse domain remain in proximity of each other.

We show that, for any level of privacy, the reported outcomes for local differential privacy (i.e., the median of the noisy time series) converges to the percentile statistics of the time series without noise as the number of time series (i.e., the size of the dataset) grows.

---

[2]While this might seem like an pathological scenario (half at minimum, half at maximum, one floating in between), privacy mechanisms must cope with all possible scenarios, not just those that are considered more realistic.

[3]An atom-less density function is a density function that assigns zero probability to any Lebesgue measure-zero event. All probability density functions with continuously-differentiable cumulative probability functions are atom-less.

This implies that, as long as there are enough data points, local differential privacy outperforms differential privacy in terms of utility. This is experimentally observed for real energy data as well if the reported percentile statistics are based on a large sample size (i.e., it is not required to report based on few outliers as in the case of 5% or 95% percentiles).

- *Stochastic Differential Privacy*: We consider a relaxation of differential privacy referred to as randomized, probabilistic, or stochastic differential privacy; see, e.g., [Machanavajjhala et al., 2008, Rubinstein and Aldà, 2017, Hall et al., 2012] for definition and applications. Stochastic differential privacy requires the definition of differential privacy to hold for *almost all* datasets (with a pre-selected high probability) rather than *all* datasets (as is the case in differential privacy). This results in ignoring cases that make the sensitivity of the reporting function to variations of individual entries in the dataset large as they occur with negligible probability if the datasets are randomly distributed and the size of the dataset is large. We show that stochastic differential privacy rivals the performance of the local differential privacy; however, it provides a weaker guarantee as the definition is relaxed to almost all cases rather then all.

- *Information-Theoretic Privacy*: Finally, we use the entropy of the additive noise as a measure of information leakage for developing an information-theoretic notion of privacy. The choice of the entropy is motivated by the use of Fano's inequality in information theory [Fano, 1961], which relates the ability of a colluding adversary (i.e., an adversary that has access to the information of everyone within the society except one individual) for estimating the data of the individual whose private information is missing to the conditional entropy of the private information. The conditional entropy of the private information is lower bounded by the entropy of the additive noise. We obtain optimal privacy-preserving policies by maximizing the entropy of the additive noise subject to various notions of utility:

    - *Expectation of the Absolute Value*: We assume that the utility is inversely proportional to the expectation of absolute value of the additive noise. Therefore, if expectation of the absolute value of the additive noise is smaller than an *a priori* bound, the utility will be above a preferred level. We show that the optimal additive noise, in this case, is Laplace, which would also guarantee differential privacy. This enables us to relate differential privacy and information-theoretic privacy.

    - *Variance*: In this case, the utility is inversely proportional to the variance of the noise. We show that the optimal additive noise is Gaussian.

    - *Essential Extrema*: Bounds on the expectation of the absolute value and the variance allow for somewhat large corruptions, albeit with relatively small probabilities. We assume that the utility is inversely proportional to the essential extrema of the noise (the extrema in the absence of a measure zero set). By ensuring that the additive noise is within a bounded set, the utility will be above a preferred level. We show that the optimal additive noise is uniformly distributed.

Generally, the guarantees of the information-theoretic privacy are weaker because they

require assumptions on the underlying dataset (e.g., that it is statistically distributed or that it complies to a certain prior distribution), the nature of the adversary must be well-understood (e.g., a colluding adversary in this report), and the side channel information to the adversary must be modelled (in this report, it is only the information of all other colluding data owners and no other additional outside source). However, due to directly trading-off privacy and utility in an optimization problem, we can understand the implications of privacy on the usefulness of the reported outcomes easily. This is contrary to differential privacy that requires numerical investigation of the utility for the specific dataset.

## 1.2 Privacy-Preserving Clustering

In Chapter 3, we investigate the privacy risks of clustering household consumption data as well as the use of local differential privacy for extracting privacy-preserving policies.

We follow the clustering algorithm in [Motlagh et al., 2019]. The proposed method relies on delay coordinate embedding for the clustering of household consumption data. This is a model-based approach which can handle unequal, asynchronous, and noisy time series. In order to understand privacy risks of clustering, we investigate removing an individual household from the clustering dataset and compute the clusters again. If the clusters are privacy-preserving, the new clusters do not change in comparison to the old clusters. We use the open dataset of smart meter measurements [Department of the Environment and Energy, 2015], also used in [Motlagh et al., 2019], for evaluation of the clustering algorithms and their privacy risks. We particularly use the data of 10,905 homes with at least one full year of time series data. This dataset is referred to as the Smart-Grid Smart-City (SGSC) dataset. Noting that the clusters change massively, up to 2% of the households change clusters, there might be a risk of privacy infringement from releasing the clusters without privacy treatments. The changes in the clustering can also be caused by the underlying clustering algorithm being robust to changes, i.e., the algorithm might only manage to find a locally optimal cluster from many local solutions and, with a small change, the algorithm might recover another locally optimal solution.

To alleviate the risk, we propose the use of local differential privacy. Since the clustering is performed on a delay coordinate embedding of the time series, instead of adding noise to the original load consumptions, we add noise the to the parameters of the delay coordinate embedding of the time series. We use Jaccard metric to evaluate the utility loss by matching the corresponding clusters, i.e. the clusters without noise and clusters with privacy-preserving noise. We compute an average Jaccard similarity index between the clusters with and without privacy-preserving noise. The differential privacy noise changes the clusters massively, in terms of the average Jaccard similarity, for reasonable choices of differential privacy parameters. However, as the differential privacy budget parameter increases, i.e., the privacy guarantee weakens, the clusters become more similar. This further points to the sensitive nature of the clusters and their possible privacy risks.

# 2    Privacy-Preserving Release of Percentile Statistics of Energy Data

In this chapter, we consider the problem of privacy-preserving release of percentile statistics of energy data. In Section 2.1, we develop privacy-preserving policies for the release of percentile statistics of energy data using differential privacy, local differential privacy, and stochastic differential privacy. In Section 2.2, we present an information-theoretic platform for privacy-aware release of energy data statistics. Section 2.1.4 presents numerical results based on a public dataset of half-hour energy measurements from 300 homes in Australia with rooftop solar systems[1].

## 2.1    Differential Privacy and Local Differential Privacy

Consider a private dataset of the form $(x_i(k))_{k\in\mathbb{K},i\in\mathbb{I}}$ with $\mathbb{K} := \{0,\ldots,K-1\}$ and $\mathbb{I} := \{1,\ldots,n\}$ in which $K \in \mathbb{N}$ (where $\mathbb{N}$ is the set of natural numbers) denotes the time horizon and $n \in \mathbb{N}$ denotes the number of the private entries in the dataset. For instance, $x_i(k)$ may denote the consumption of household $i \in \mathbb{I}$ at time instant $k \in \mathbb{K}$ measured by a smart meter. It is desired to report a time-series signal of the form $f((x_i(k))_{i\in\mathbb{I}})$, $\forall k \in \mathbb{K}$, where $f(\cdot)$ is a percentile statistic, e.g., the $50\%$ percentile statistic or, as more commonly known, the median. Assume that, for the sake of preserving the privacy of participants, the following output is instead released:

$$z(k) = f((x_i(k))_{i\in\mathbb{I}}) + w(k), \quad \forall k \in \mathbb{K}, \tag{2.1}$$

where $w := (w(k))_{k\in\mathbb{K}}$ is a random variable with density function $\gamma$. Throughout the paper, for any time series $(\xi(k))_{k\in\mathbb{K}}$, define $\xi := (\xi(k))_{k\in\mathbb{K}}$. The density of the additive noise can be selected to ensure differential privacy, formally defined below.

**Definition 2.1.1** (Differential Privacy). *The additive mechanism* (2.1) *is $\epsilon$-differentially private if, for all Borel measurable sets $\mathbb{B} \subseteq \mathbb{R}^K$,*

$$\mathbb{P}\{z \in \mathbb{B} \,|\, (x_i)_{i\in\mathbb{I}}\} \leq \exp(\epsilon)\mathbb{P}\{z \in \mathbb{B} \,|\, (\bar{x}_i)_{i\in\mathbb{I}}\}, \ \forall(x,\bar{x}) \in \mathbb{A},$$

*where $((x_i)_{i\in\mathbb{I}}, (\bar{x}_i)_{i\in\mathbb{I}}) \in \mathbb{A}$ states that datasets $(x_i)_{i\in\mathbb{I}}$ and $(\bar{x}_i)_{i\in\mathbb{I}}$ are adjacent.*

In the remainder of this section, various notions of adjacency $\mathbb{A}$ are explored. However, before that, a few other useful concepts are presented. Note that, for any notion of adjacency, the following sensitivity can be computed:

$$\Delta f = \max_{((x_i)_{i\in\mathbb{I}}, (\bar{x}_i)_{i\in\mathbb{I}})\in\mathbb{A}} \sum_{k\in\mathbb{K}} |f((x_i(k))_{i\in\mathbb{I}}) - f((\bar{x}_i(k))_{i\in\mathbb{I}})|.$$

The following theorem proves that adding Laplace noise with an appropriate scaling parameter depending on the above-mentioned sensitivity provides differential privacy.

---

[1]https://www.ausgrid.com.au/Common/About-us/Corporate-information/Data-to-share/Solar-home-electricity-data.aspx

**Theorem 2.1.1** ([Dwork and Roth, 2014])**.** *The additive mechanism* (2.1) *is $\epsilon$-differentially private if* $(w(k))_{k \in \mathbb{K}}$ *are i.i.d.[2] Laplace noises with scale $\zeta/\epsilon$ for any constant $\zeta \geq \Delta f$.*

An alternative to adding noise after calculating the percentile statistic $f(\cdot)$ is to add noise to the raw data first and then compute the percentile statistics as

$$z(k) = f((y_i(k))_{i \in \mathbb{I}}), \qquad\qquad \forall k \in \mathbb{K}, \qquad\qquad (2.2a)$$

$$y_i(k) = x_i(k) + v_i(k), \qquad\qquad \forall k \in \mathbb{K}, \qquad\qquad (2.2b)$$

where $v_i := (v_i(k))_{k \in \mathbb{K}}$ is a random variable. This approach is of particular interest when the data aggregator, e.g., electricity retailer, is untrusted because each household can add an appropriate noise to guarantee differential privacy at an individual level, referred to as local differential privacy.

**Definition 2.1.2** (Local Differential Privacy)**.** *The additive mechanism* (2.2) *is local $\epsilon$-differentially private if, for all Borel measurable sets* $\mathbb{B} \subseteq \mathbb{R}^K$,

$$\mathbb{P}\{y_i \in \mathbb{B} \mid x_i\} \leq \exp(\epsilon)\mathbb{P}\{y_i \in \mathbb{B} \mid \bar{x}_i\}, \forall(x_i, \bar{x}_i) \in \mathbb{A}_i, \forall i \in \mathbb{I},$$

*where* $(x_i, \bar{x}_i) \in \mathbb{A}_i$ *states that time series* $x_i$ *and* $\bar{x}_i$ *are (locally) adjacent.*

To investigate the reporting mechanism (2.2), define

$$\Delta X = \max_{((x_i)_{i \in \mathbb{I}}, (\bar{x}_i)_{i \in \mathbb{I}}) \in \mathbb{A}} \max_{i \in \mathbb{N}} \|x_i - \bar{x}_i\|_1.$$

The following theorem shows that adding Laplace noise before computing the percentile statistics also provides local differential privacy.

**Theorem 2.1.2.** *The mechanism* (2.2) *is local $\epsilon$-differentially private if* $(v_i(k))_{k \in \mathbb{K}}$ *are i.i.d. Laplace noises with scale $\zeta/\epsilon$ for any constant $\zeta \geq \Delta X$.*

*Proof.* The proof follows from the same line of reasoning as in Theorem 2.1.1. □

**Remark 2.1.1** (Interpretation of $\epsilon$ in Differential Privacy)**.** *An adversary with access to* $z(k)$ *wants to use hypothesis testing to detect a characteristic or trait of household $i$ (e.g., if the house is empty or not). Then, the Chernoff-Stein Lemma states that the logarithm of the smallest probability of miss-detection (i.e., the event that null hypothesis is accepted while the alternative hypothesis is true) for a fixed probability of false negative (i.e., the event that null hypothesis is rejected while it is true) scales negatively with Kullback-Leibler divergence of the probability density of* $z(k)$ *if the null hypothesis holds* $p_0$ *and the probability density of* $z(k)$ *if the alternative hypothesis holds* $p_1$ *[Wang et al., 2009]. Thus, if the Kullback-Leibler divergence is small, the smallest probability of miss-detection is large and vice versa. The Kullback-Leibler divergence of* $p_0$ *and* $p_1$ *is given by*

$$\int p_0(z) \log\left(\frac{p_0(z)}{p_1(z)}\right) dz \leq \int p_0(z) \log(\exp(\epsilon)) dz = \epsilon,$$

*where the inequality follows from that* $p_0(z)/p_1(z) \leq \exp(\epsilon)$ *by the definition of differential privacy, if the additive noise is atom-less (i.e., it admits a probability density function), e.g., the*

---

[2]i.i.d. stands for identically and independently distributed.

*Laplace additive noise in Theorem 2.1.1. Therefore, reducing $\epsilon$ makes the Kullback-Leibler divergence smaller and that implies that the smallest probability of miss-detection grows larger (rending the hypothesis test by the adversary useless). Now, assume that the adversary wants to estimate the energy time series from household $i$. In that case, the Cramér-Rao bound [Shao, 2003, p. 169] illustrates that*

$$\frac{1}{K}\sum_{k\in\mathbb{K}}\mathbb{E}\{(x_i(k)-[\hat{x}_i(z)](k))^2\}=\frac{\mathbb{E}\{\|x_i-\hat{x}_i(z)\|_2^2\}}{K}\geq\frac{\Delta f^2}{\epsilon^2},$$

*where $\hat{x}_i(z)$ denotes any unbiased estimate of $x_i$ based on the privacy-preserving outputs $z$ when using the additive noise in Theorem 2.1.1. Clearly, by reducing $\epsilon$, the estimation error grows.*

In the following subsections, various notions of adjacency are investigated. Capturing the "smallest" adjacency set is most beneficial as it effectively reduces $\Delta f$ and $\triangle X$, which in turn makes the scale of the additive noise smaller; see Theorem 2.1.1 and 2.1.2. The adjacency set should contain all the possibilities for time-series extracted from smart meters as otherwise the privacy guarantee is weakened (i.e., time series that are not captured can potentially have significant effects on the reported outcome and thus be identifiable).

### 2.1.1 Point-wise Adjacency

We start by considering a simple concept of adjacency defined below.

**Assumption 2.1.1** (Boundedness)**.** *There exists $x_{\max} > 0$ such that, for all $i$, $x_i \in \mathbb{X}$, where $\mathbb{X} = [-x_{\max}, x_{\max}]^K$.*

**Definition 2.1.3** (Point-wise Adjacency)**.** *The point-wise adjacency is defined as*
$\mathbb{A} := \{((x_i)_{i\in\mathbb{I}}, (\bar{x}_i)_{i\in\mathbb{I}}) \in \mathbb{X}^n \times \mathbb{X}^n \,|\, x_j(k) \neq \bar{x}_j(k)$ *for at most one $j \in \mathbb{I}$ and $k \in \mathbb{K}\}$.*

Definition 2.1.3 states that two datasets $(x_i)_{i\in\mathbb{I}}$ and $(\bar{x}_i)_{i\in\mathbb{I}}$ are adjacent, or $((x_i)_{i\in\mathbb{I}}, (\bar{x}_i)_{i\in\mathbb{I}}) \in \mathbb{A}$, if they differ from each other in at most one time instant and for at most one time series (i.e., one consumption profile). For point-wise adjacency, we have $\Delta f = \Delta X = 2x_{\max}$.

**Corollary 2.1.1.** *For the point-wise adjacency in Definition 2.1.3, the additive mechanism* (2.1) *is $\epsilon$-differentially private if $(w(k))_{k\in\mathbb{K}}$ are i.i.d. Laplace noises with scale $2x_{\max}/\epsilon$.*

Point-wise adjacency, and in turn Corollary 2.1.1, does not assume anything more than that the dataset lies within $[-x_{\max}, x_{\max}]^K$, which can be potentially very large. This adjacency ignores the fact that the dataset is a time-series in which the consumption at a given time instant is somewhat correlated to the previous time instances.

**Definition 2.1.4** (Local Point-wise Adjacency)**.** *The local point-wise adjacency is defined as*
$\mathbb{A}_i := \{(x_i, \bar{x}_i) \in \mathbb{X} \times \mathbb{X} \,|\, x_j(k) \neq \bar{x}_j(k)$ *for at most one $k \in \mathbb{K}\}$.*

Note that $((x_i)_{i\in\mathbb{I}}, (\bar{x}_i)_{i\in\mathbb{I}}) \in \mathbb{A}$ if, for at most one $j \in \mathbb{I}$, $(x_j, \bar{x}_j) \in \mathbb{A}_j$.

**Corollary 2.1.2.** *For the local point-wise adjacency in Definition 2.1.4, the mechanism* (2.2) *is $\epsilon$-differentially private if $(v(k))_{k\in\mathbb{K}}$ are i.i.d. Laplace noises with scale $2x_{\max}/\epsilon$.*

Corollary 2.1.2 states the same level of noise that guarantees differential privacy for percentile statistics, Corollary 2.1.1, if first added to the raw data can also guarantees local differential

privacy. This is due to the very nature of percentile statistics. One can always construct worst-case scenarios in which arbitrary variations of a single entry can cause arbitrary variations of the percentile statistics. This is unfortunately a drawback of using differential privacy when dealing with percentile statistics. The benefit of adding the noise to the raw data versus the percentile statistics is that the outcome is statistically reasonable (i.e., percentile statistics are sorted and the outcome may not require post processing if $n$ is large enough) and the use of percentile statistics can reduce the error in the reported percentile statistics. There always exists $N \in \mathbb{N}$ such that

$$\mathbb{E}\{|f((y_i(k))_{i\in\mathbb{I}}) - f((x_i(k))_{i\in\mathbb{I}})|\} < \mathbb{E}\{|z(k) - f((x_i(k))_{i\in\mathbb{I}})|\}, \; \forall n \geq N,$$

where $w_i(k)$ and $v_i(k)$ are given by Corollaries 2.1.1 and 2.1.2. In fact, it can be proved that $\lim_{n\to\infty} \mathbb{E}\{|f((y_i(k))_{i\in\mathbb{I}}) - f((x_i(k))_{i\in\mathbb{I}})|\} = 0$.

The following definition relaxes the differential privacy condition to hold for almost all datasets (in a statistical sense). Such a relaxations allows us to avoid pathological cases that make $\Delta f$, and the scale of the additive noise, large. This notion of privacy has been previously proposed and investigated in [Machanavajjhala et al., 2008, Rubinstein and Aldà, 2017, Hall et al., 2012] (under the names of random differential privacy or probabilistic differential privacy). The relaxation comes at the price of the assumption that the dataset is statistically distributed, which is not required in differential privacy.

**Definition 2.1.5** (Stochastic Differential Privacy). *Assume that $x_i := (x_i(k))_{k\in\mathbb{K}}$ are i.i.d. random variables. The additive mechanism* (2.1) *is $\epsilon$-differentially private with probability, at least, $1 - \varrho$ if, for all Borel measurable sets $\mathbb{B} \subseteq \mathbb{R}^K$,*

$$\mathbb{P}\left\{\mathbb{P}\{z \in \mathbb{B}|(x_i)_{i\in\mathbb{I}}\} \leq \exp(\epsilon)\mathbb{P}\{z \in \mathbb{B}|(\bar{x}_i)_{i\in\mathbb{I}}\}, \forall(x, \bar{x}) \in \mathbb{A}\right\} \geq 1 - \varrho.$$

Define $p_{\min} = \min_{k\in\mathbb{K}} p_k(\nu(k))$, $p_k(x(k)) = \int p(x(k), (x(\ell))_{\ell\in\mathbb{K}\setminus\{k\}})\mathrm{d}\mu((x(\ell))_{\ell\in\mathbb{K}\setminus\{k\}})$, and $\nu[k] \in \{x(k) \mid \int_{-\infty}^{x(k)} p_k(\alpha)\mathrm{d}\alpha = 1/2\}$.

**Corollary 2.1.3.** *Assume that $x_i := (x_i(k))_{k\in\mathbb{K}}$ are i.i.d. samples from a probability density function $p$ with a lower bound away from zero. For large enough $n$, using the point-wise adjacency in Definition 2.1.3, and for $f(\cdot) = \mathrm{median}(\cdot)$, the additive mechanism* (2.1) *is $\epsilon$-differentially private with probability, at least, $1 - \varrho$ if $(w(k))_{k\in\mathbb{K}}$ are i.i.d. Laplace noises with scale $\min\{((2n-1)/(n(n-1)p_{\min}^2\varrho))^{1/2}, 2x_{\max}\}/\epsilon$.*

Corollary 2.1.3 shows that, for large $n$, differential privacy can be guaranteed with a high probability using a negligible additive noises. In fact, $\lim_{n\to 0} \min\{((2n-1)/(n(n-1)p_{\min}^2\varrho))^{1/2}, 2x_{\max}\}/\epsilon = 0$ for all $\varrho > 0$ Therefore, as $n$ goes to infinity, the mechanism with no additive noise becomes differentially private with probability one! Note that, in Corollary 2.1.3, $p_{\min}$ is the value of the density function evaluated at the median of the density function. This can be easily proved for all other percentile statistics.

### 2.1.2 Trajectory Adjacency

Time series are often outcomes of physical systems pointing to their smoothness and therefore, alternative notions of adjacency might be more suitable for them.

**Definition 2.1.6** (Trajectory Adjacency)**.** *The local trajectory adjacency is defined as* $\mathbb{A}_i := \&(x_i, \bar{x}_i) \in \mathbb{X} \times \mathbb{X} \,|\, |x_j(k) - \bar{x}_j(k)| \leq \rho, \forall k \in \mathbb{K}\}$. *The trajectory adjacency* $\mathbb{A}$ *is defined as the set of all* $((x_i)_{i\in\mathbb{I}}, (\bar{x}_i)_{i\in\mathbb{I}})$ *such that, for at most one* $j \in \mathbb{I}$, $(x_j, \bar{x}_j) \in \mathbb{A}_j$ *and* $x_i(k) = \bar{x}_i(k)$, $\forall k \in \mathbb{K}$, *for* $i \neq j$.

Definition 2.1.3 states that two datasets $(x_i)_{i\in\mathbb{I}}$ and $(\bar{x}_i)_{i\in\mathbb{I}}$ are adjacent, or $((x_i)_{i\in\mathbb{I}}, (\bar{x}_i)_{i\in\mathbb{I}}) \in \mathbb{A}$, if they differ from each other for at most one time series and, for the differing time series, the values remain within a tube of the size of $\rho$ of each other. Similarly, for the trajectory adjacency in Definition 2.1.6, it can be proved that $\Delta f = \Delta X = 2\rho K$.

**Corollary 2.1.4.** *For the trajectory adjacency in Definition 2.1.6, the additive mechanism* (2.1) *is* $\epsilon$-*differentially private if* $(w(k))_{k\in\mathbb{K}}$ *are i.i.d. Laplace noises with scale* $2\rho K/\epsilon$.

Depending on the value of $\rho$ this can be more relaxed in comparison to the point-wise adjacency case; however, meeting differential privacy becomes harder as $K$ grows.

**Corollary 2.1.5.** *For the local trajectory adjacency in Definition 2.1.6, the mechanism* (2.2) *is local* $\epsilon$-*differentially private if* $(v_i(k))_{k\in\mathbb{K}}$ *are i.i.d. Laplace noises with scale* $2\rho K/\epsilon$.

### 2.1.3    Sparse Domain Adjacency

Time series associated with physical systems are most often sparse in another domain, e.g., the frequency domain. This prompt us to work with the underlying database within a sparsifying domain. Let $\Omega \in \mathbb{R}^{K \times K}$ be a dictionary that sparsifies the dataset. The equivalent dataset in the other domain is given by $(X_i)_{i\in\mathbb{I}} := (\Omega x_i)_{i\in\mathbb{I}}$.

**Remark 2.1.2** (Discrete Cosine Transform)**.** *For instance, we can represent the data in the frequency domain using the Discrete Cosine Transform (DCT). Define the matrix* $\Omega \in \mathbb{R}^{K \times K}$ *such that the entry on the* $i$-*th row and the* $j$-*th column is given by*

$$\omega_{ij} := \sqrt{2/K} \cos((\pi/K)(j - 1/2)(i - 1/2)).$$

*The frequency domain equivalent of a time series* $x = (x(k))_{k\in\mathbb{K}}$ *is given by* $X = \Omega x$. *Note that, by construction* $\Omega^2 = I$ *and, therefore,* $x = \Omega X$.

**Definition 2.1.7** (Sparsity)**.** *A time-series* $x_i$ *is* $\sigma$-*approximately* $L$-*sparse if* $\mathrm{card}(\{k \in \mathbb{K} \,|\, |[\Omega x_i](k)| \geq \sigma\}) \leq L$. *The set of all* $\sigma$-*approximately* $L$-*sparse time-series is denoted by* $\mathbb{L}(L, \sigma)$. *Further, the time-series is (exactly)* $L$-*sparse if it is* $0$-*approximately* $L$-*sparse.*

A $\sigma$-approximately $L$-sparse sparse time-series $x_i$ can be perturbed to get another $\sigma$-approximately $L$-sparse time-series $\tilde{x}_i$ as

$$\Omega(\Omega x_i + \bar{v}_i) = x_i + \Omega \bar{v}_i = x_i + v_i, \tag{2.3}$$

where $v_i := \Omega \bar{v}_i$ and $\bar{v}_i$ is a vector of noise such that $\bar{v}_i(k) = 0$ if $k \notin \{t \in \mathbb{K} \,|\, |x_i(t)| \geq \sigma\}$ and, otherwise, $\bar{v}_i(k)$ is an i.i.d. Laplace noise with scale $b$ (or variance $2b^2$).

**Definition 2.1.8** (Local Sparse-Domain Adjacency)**.** *The local sparse-domain adjacency is defined as* $\mathbb{A}_i := \{(x_i, \bar{x}_i) \in \mathbb{L}(L, \sigma) \times \mathbb{L}(L, \sigma) \,|\, |[\Omega x_i](k) - [\Omega \bar{x}_i](k)| \leq \varpi, \forall k \in \mathbb{K}\}$.
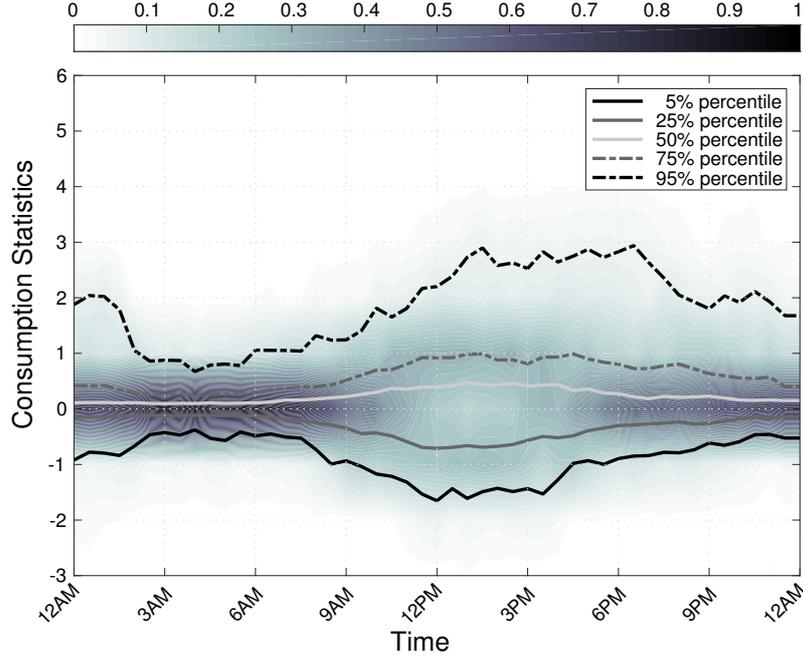
**Figure 2.1: The probability density function of the electricity consumption of the households in the dataset and the corresponding percentile statistics without implementing a privacy-preserving policy. The curves illustrate the 5% (——), 25% (——), 50% (——), 75% (- - -), and 95% (- · -) percentiles and the intensity captures probability density function of the consumption.**

**Corollary 2.1.6.** *The perturbed mechanism in* (2.3) *with Laplace noises with scale* $\zeta/\epsilon$ *is locally* $\epsilon$*-differentially private for the family of* $\sigma$*-approximately* $L$*-sparse datasets,* $\mathbb{L}(L, \sigma)$*, with any constant* $\zeta \geq 2L\varpi + 2(K - L)\sigma$.

### 2.1.4 Numerical Analysis

We start by presenting the statistics of the dataset in the *absence of a privacy-preserving policy*. Figure 2.1 illustrates the probability density function of the energy consumptions of the households in the dataset and the corresponding percentile statistics. The solid curves illustrate the 5%, 25%, 50%, 75%, and 95% percentiles. As somewhat easy to see from the figure, $x_{\max} = 4$.

Noting that the scaling parameter of the additive Laplace noise must be set as $2x_{\max}/\epsilon = 8/\epsilon$ in the case of point-wise adjacency, we can see that the magnitude of the additive noise is fairly large in comparison to the original statistics even for modest choices of $\epsilon$. The signal to noise ratio for the median is in fact equal to

$$\mathrm{SNR} := \sqrt{\frac{\sum_{k=1}^{T} f((x_i(k))_{i \in \mathbb{I}})^2}{\mathbb{E}\{\|w\|_2^2\}}} \approx 2.4 \times 10^{-2}\epsilon$$

Hence, even for moderate choices of $\epsilon$, $\mathrm{SNR}$ is abhorrent and so follows the privacy-utility trade-off.

Figure 2.2 illustrates the percentile statistics of the households with the privacy-preserving policy in Corollary 2.1.1 for $\epsilon = 20$. Note that, here, some post processing has been done to
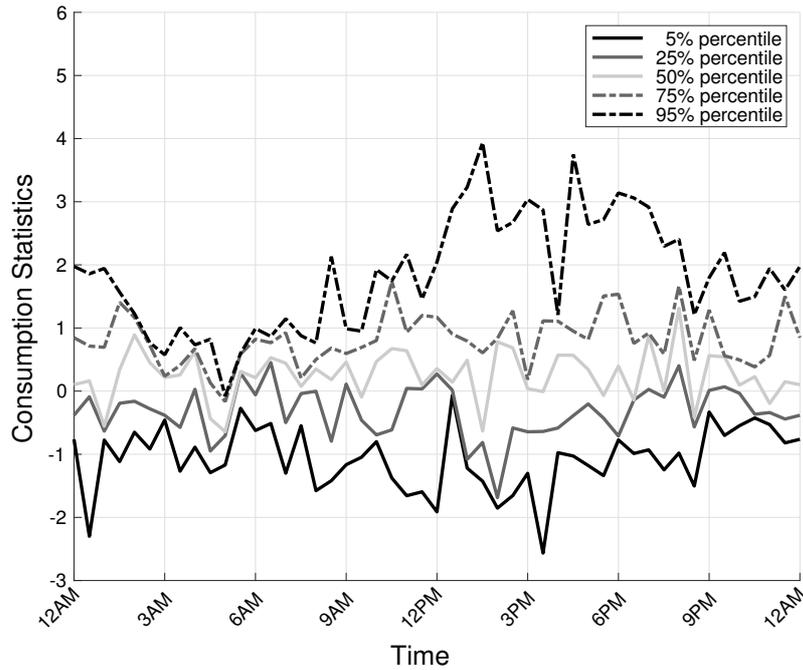
**Figure 2.2: The reported percentile statistics of the households with the privacy-preserving policy in Corollary 2.1.1 for $\epsilon = 20$.**

ensure that the percentile statistics make sense. This is required because, by adding independently and identically distributed (i.i.d.) noises to the percentile statistics, their relative order might be changed (e.g., median might become smaller than the $25\%$ percentile) which is not acceptable/plausible. Therefore, the corrupted statistics are sorted after adding the noise. Post processing does not weaken the privacy guarantee [Dwork and Roth, 2014]. For $\epsilon = 20$, $\mathbb{E}\{(z(k) - f((x_i(k))_{i\in\mathbb{I}}))^2\} = 128/\epsilon^2 = 0.32$ (which is comparable to the magnitude of the percentile statistics, such as the median). Due to the magnitude of the additive noise, the statistics are highly distorted. We should note that, because five percentile statistics are reported, the overall privacy guarantee is five times weaker (i.e., $\epsilon = 100$ for the overall case) [Dwork and Roth, 2014].

Now, we can compare the effect of differential privacy in Figure 2.2 with local differential privacy in Figure 2.3. Figure 2.3 illustrates the difference between the reported output $f((y_i(k))_{i\in\mathbb{I}})$ and the original statistics $f((x_i(k))_{i\in\mathbb{I}})$ with the privacy-preserving policy in Corollary 2.1.2 for $\epsilon = 20$. The straight solid curve with circle markers shows the $\mathbb{E}\{(z(k) - f((x_i(k))_{i\in\mathbb{I}}))^2\} = 0.32$, which is the variance of the reporting error when using the differentially private policy in Corollary 2.1.1. In the case of local differential privacy, $\mathbb{E}\{(z(k) - f((x_i(k))_{i\in\mathbb{I}}))^2\}$ is empirically computed with $10^4$ particles. It is important to the note that the effect of the additive noise (captured by the degradation in the usefulness of the data) is more drastic for the outliers. This can be seen in Figure 2.3 by that the 5% and 95% percentiles are most affected by the privacy-preserving noise. Finally, note that generating any number of percentile statistics does not reduce the privacy guarantee if all such statistics are generated with the same additive noise on the data (i.e., from the same perturbed dataset with local differential privacy); however, the privacy guarantee gets weaker if a noisy dataset is
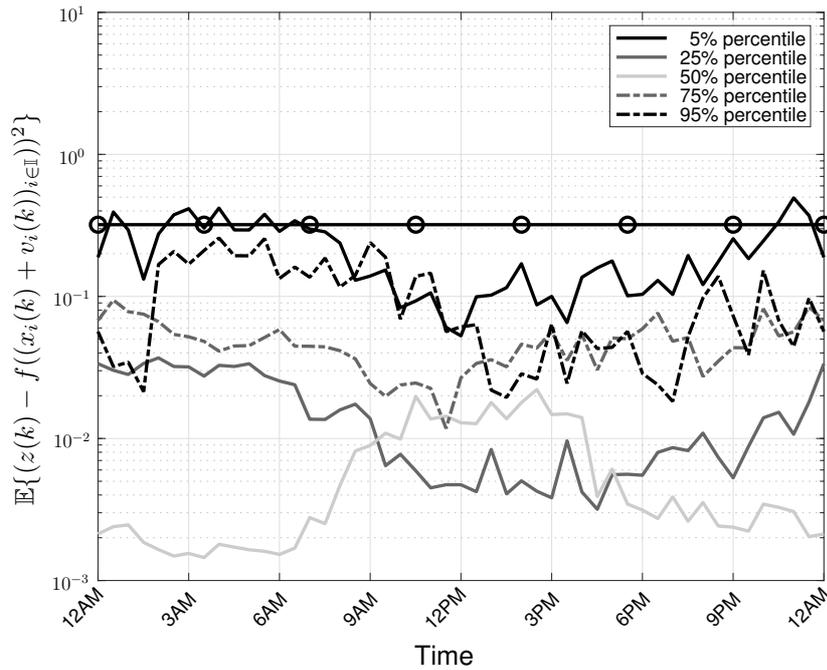
**Figure 2.3:** The difference between the reported values $f((y_i(k))_{i\in\mathbb{I}})$ and the original statistics $f((x_i(k))_{i\in\mathbb{I}})$ with the privacy-preserving policy in Corollary 2.1.2 for $\epsilon = 20$ (i.e., local differential privacy). The curve (–◦–) illustrates the difference between the reported values and original percentile statistics for the privacy-preserving policy in Corollary 2.1.1 for $\epsilon = 20$ (i.e., differential privacy). The utility of the privacy-preserving policy in Corollary 2.1.2 is better than the privacy-preserving policy in Corollary 2.1.1 for almost all percentile statistics.

generated for computing each statistics. Evidently, in Figure 2.3, the error in the median is negligible as predicted earlier.

Now, we can consider the trajectory adjacency with $\rho = 0.1$. Figure 2.4 shows the difference between the reported outputs $f((y_i(k))_{i\in\mathbb{I}})$ and the original statistics $f((x_i(k))_{i\in\mathbb{I}})$ with the privacy-preserving policy in Corollary 2.1.5 for $\epsilon = 20$. The straight solid curve with circle markers shows the $\mathbb{E}\{(z(k) - f((x_i(k))_{i\in\mathbb{I}}))^2\}$ for differentially private policy in Corollary 2.1.4. The utility-privacy is still not good for this notion of adjacency in the case of differential privacy; however, as expected, the privacy-utility trade-off of local differential privacy is acceptable.

We can now demonstrate the sparse adjacency using DCT transform. Figure 2.5 shows the difference between the reported outputs $f((y_i(k))_{i\in\mathbb{I}})$ and the original statistics $f((x_i(k))_{i\in\mathbb{I}})$ with the privacy-preserving policy in Proposition 2.1.6 for $\epsilon = 20$ and assuming that the underlying dataset is 5-sparse and $\varpi = 0.5$. In practice, the dataset might not be exactly 5-sparse, which might make the privacy guarantee weaker; however, the utility-privacy trade-off is as good as the case of local differential privacy for both point-wise and trajectory adjacencies. To validate the assumption of 5-sparsity, we can check the frequency content of the energy data. Figure 2.6 shows the frequency content. In this figure, the notation $\sigma_j(\Omega x_i)$ denotes the $j$-th largest component of $\Omega x_i$. The first five components (i.e., the components that are equal to or larger than the solid black horizontal line) contain 88% of the energy (note that the energy is proportional to the square value of the component). Therefore, the 5-sparsity assumption is in line with the nature of the dataset.
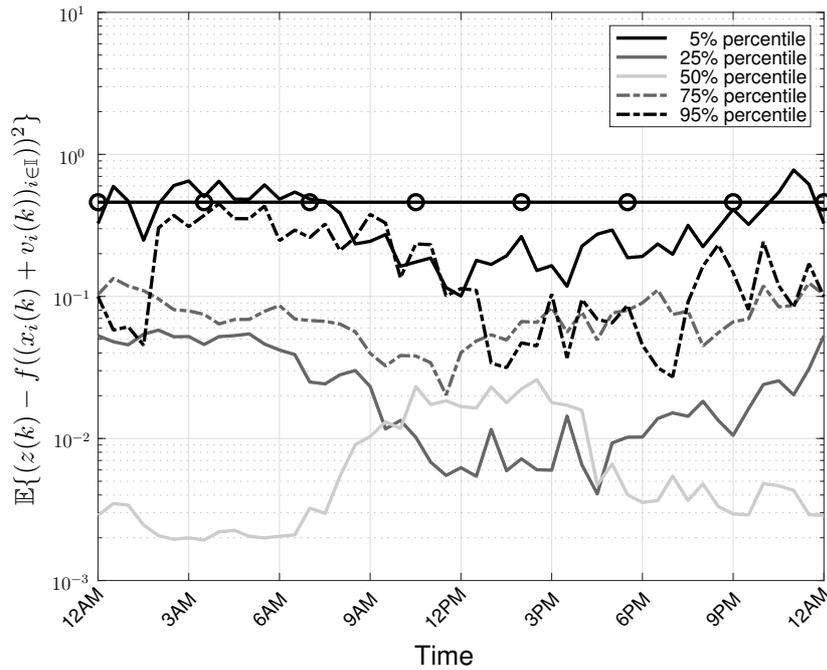
**Figure 2.4: The difference between the reported values** $f((y_i(k))_{i \in \mathbb{I}})$ **and the original statistics** $f((x_i(k))_{i \in \mathbb{I}})$ **with the privacy-preserving policy in Corollary 2.1.5 for** $\epsilon = 20$ **assuming that** $\rho = 0.1$ **(i.e., local differential privacy). The curve (—o—) illustrates the difference between the reported values and original percentile statistics for the privacy-preserving policy in Corollary 2.1.4 for** $\epsilon = 20$ **assuming that** $\rho = 0.1$ **(i.e., differential privacy). The utility of the privacy-preserving policy in Corollary 2.1.5 is better than the privacy-preserving policy in Corollary 2.1.4 for almost all percentile statistics.**

Table 2.1 summarises all the points discussed in this subsection for the numerical example.

## 2.2    Information-Theoretic Analysis

This section addresses privacy from an information-theoretic perspective. The following assumption is made through out the information-theoretic analysis.

**Assumption 2.2.1.** $x_i := (x_i(k))_{k \in \mathbb{K}}$ *are i.i.d. samples from a probability density function* $p$.

Assumption 2.2.1 is not very conservative as otherwise reporting statistics, such as median, is meaningless. Note that knowledge of the distribution *a priori* is not required for the following analysis.

**Table 2.1: Average squared error between the reported output and the original statistics (averaged over the horizon) for point-wise, trajectory, and sparse adjacency with** $\epsilon = 20$**. Local differential privacy clearly outperforms differential privacy in all cases.**

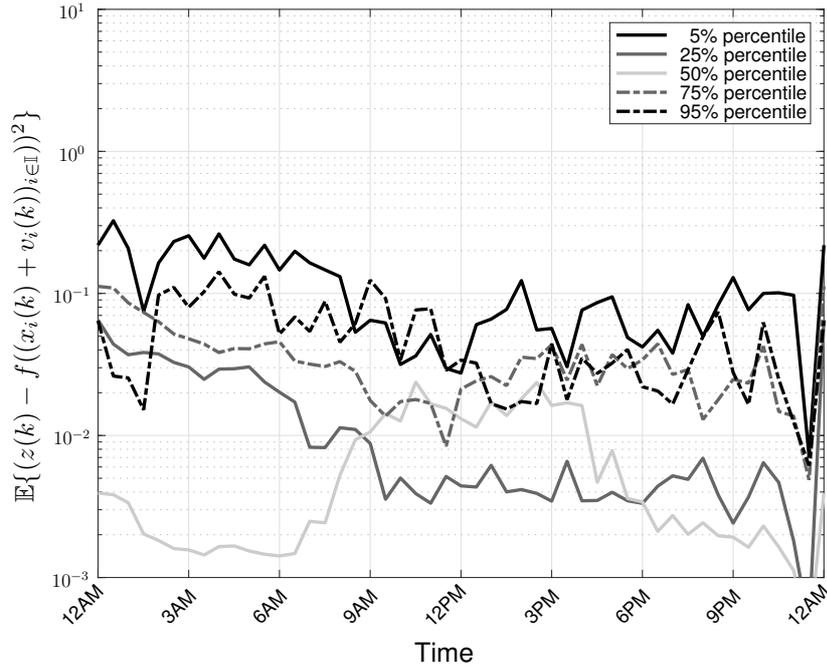|  | Point-wise Adjacency | | Trajectory Adjacency | | Sparse Adjacency |
|---|---|---|---|---|---|
|  |  | Local |  | Local | Local |
|  | Differential Privacy | Differential Privacy | Differential Privacy | Differential Privacy | Differential Privacy |
| 5% percentile | 0.3200 | 0.2088 | 0.4608 | 0.3473 | 0.1116 |
| 25% percentile | 0.3200 | 0.0153 | 0.4608 | 0.0252 | 0.0142 |
| 50% percentile (median) | 0.3200 | 0.0064 | 0.4608 | 0.0079 | 0.0068 |
| 75% percentile | 0.3200 | 0.0478 | 0.4608 | 0.0723 | 0.0370 |
| 95% percentile | 0.3200 | 0.1004 | 0.4608 | 0.1751 | 0.0519 |

**Figure 2.5: The difference between the reported values** $f((y_i(k))_{i\in\mathbb{I}})$ **and original statistics** $f((x_i(k))_{i\in\mathbb{I}})$ **with the privacy-preserving policy in Proposition 2.1.6 for** $\epsilon = 20$ **and assuming that the underlying dataset is** $5$**-sparse and** $\varpi = 2$**.**

**Theorem 2.2.1** ([Cover and Thomas, 2012])**.** *Consider an all powerful adversary*[3] *who wants to extract the data of individual* $j \in \mathbb{I}$ *given the reported vector of statistics* $z$ *and the information of all the other entries of the database* $x_{-j} := (x_i)_{i\in\mathbb{I}\setminus\{j\}}$ *and denote its estimate of* $x_j$ *based on all the available information by* $\hat{x}_j(z, x_{-j})$*. Then*

$$\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_2^2\} \geq \frac{K}{2\pi e} \exp\left(\frac{2H(w)}{K}\right).$$

To ensure that $\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_2^2\}$ is large, $H(w)$ can be maximized. This leads to finding the most privacy preserving policy. On the hand, a measure of quality must be used. In the absence of a measure of quality, the most privacy preserving policy is one with an infinite entropy (essentially burying the data under noise with infinite variance).

### 2.2.1 Expected Absolute Value

One measure of privacy is the expected absolute value (or expected deviations around zero) given by $\mathbb{E}\{\|w\|_1\}$, where $\|w\|_1 = \sum_{k\in\mathbb{K}} |w(k)|$. Using this measure of quality, the following optimization problem can be solved to achieve balance between privacy and utility:

$$\mathbf{P}_1 : \max_{\gamma \in \Delta(\mathbb{R}^K) \cap \mathcal{C}^1(\mathbb{R}^K, \mathbb{R})} \quad H(w), \tag{2.4a}$$

$$\text{s.t.} \quad \mathbb{E}\{\|w\|_1\} \leq \vartheta_1 K, \tag{2.4b}$$

---

[3]This case can be seen as an extreme case in which everyone within a society wants to retrieve private data of an individual. If a reasonable guarantee in this case can be provide without sacrificing quality hugely, simpler, less sophisticated, attacks are also covered.
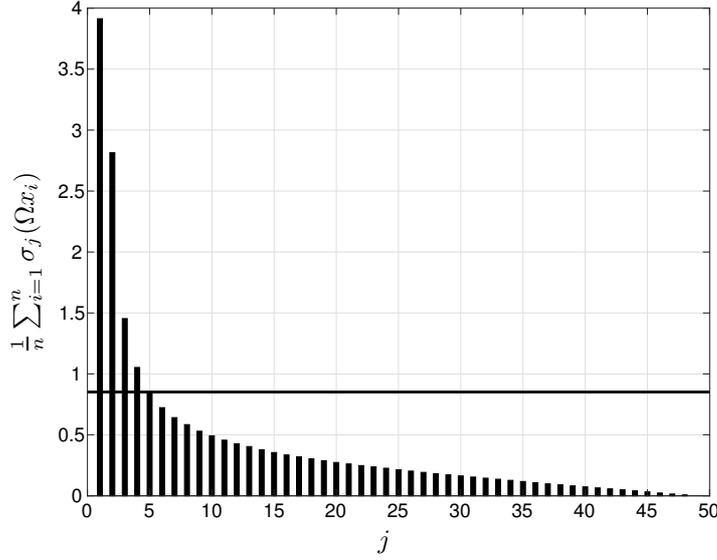
**Figure 2.6: The frequency content of the energy data using DCT transform. Here,** $\sigma_j(\Omega x_i)$ **denotes the** $j$**-th largest component of** $\Omega x_i$**. The first five components (i.e., the ones that are equal or larger than the solid black horizontal line) contain 88% of the energy of the signal.**

where $\Delta(\mathbb{R}^K)$ is the set of all probability density functions defined over $\mathbb{R}^K$ and $\mathcal{C}^1(\mathbb{R}^K, \mathbb{R})$ is the set of continuously differentiable functions from $\mathbb{R}^K$ to $\mathbb{R}$.

**Theorem 2.2.2.** *The unique solution to* $\mathbf{P}_1$ *in (2.4) is given by* $\gamma_1^*(w) = (2\vartheta_1)^{-K} \exp(-\|w\|_1/\vartheta_1)$.

Theorem 2.2.2 states that $w(k)$, $\forall k \in \mathbb{K}$, are i.i.d. Laplace random variables with zero mean and scale $\vartheta_1$ (i.e., variance of $2\vartheta_1^2$). For this policy, the privacy guarantee is given by $\mathfrak{P}(\gamma_1^*) = K/(2\pi e) \exp((2/K)H(w)) = 2e\vartheta_1^2 K/\pi$. Let us define the utility $\mathfrak{U}(\cdot)$ for a distribution to be the inverse of the quality of response. In this case, by construction, the utility is given by $\mathfrak{U}(\gamma_1^*) = 1/\mathbb{E}\{\|w\|_1\} = 1/(K\vartheta_1)$. Therefore, for the optimal policy in Theorem 2.2.2, it can be deduced that $\mathfrak{P}(\gamma_1^*)\mathfrak{U}(\gamma_1^*) = (2e\vartheta_1)/\pi = \text{const}$. This illustrates the conflict between privacy and utility, i.e., improvements in privacy always come at the cost of decreased utility.

**Corollary 2.2.1.** *The optimal privacy-preserving policy* $\gamma_1^*$ *in Theorem 2.2.2 is* $(\Delta f/\vartheta_1)$*-differential privacy.*

**Corollary 2.2.2.** *For the* $\epsilon$*-differentially private mechanism described in Theorem 2.1.1,* $(1/K)\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_2^2\} \geq 2/(\pi e)(\Delta f/\epsilon)^2$.

Corollary 2.2.2 shows that, from the perspective of the estimation error $\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_2^2\}$, $\epsilon$ does not need to stay fixed for various horizon lengths $K$ in the case of trajectory adjacency in Definition 2.1.6. For instance, if we set $\epsilon = K\epsilon_0$, it can be deduced that $(1/K)\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_2^2\} = 8\rho^2/(\pi e\epsilon_0^2)$. Therefore, to ensure that the average error $(1/K)\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_2^2\}$ is larger than or equal to $\sigma$, it must be ensure that $\epsilon_0 = \rho\sqrt{8/(\pi e\sigma)}$. If an inequality with 2-norm is not desired, one can show that $\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_\infty\} \geq ((1/K)\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_2^2\})^{1/2} = \sqrt{\sigma}$, if $\epsilon_0 = \rho\sqrt{8/(\pi e\sigma)}$. Here, $\|\xi\|_\infty = \max_{k \in \mathbb{K}} |\xi(k)|$ for any vector $\xi = (\xi(k))_{k \in \mathbb{K}}$.

### 2.2.2　Covariance

Another way of measuring the quality of the response is to use variance. In this case, we can solve the following optimisation problem to find the optimal privacy-preserving policy:

$$\mathbf{P}_2 : \max_{\gamma \in \Delta(\mathbb{R}^K) \cap \mathcal{C}^1(\mathbb{R}^K, \mathbb{R})} \quad H(w), \tag{2.5a}$$

$$\text{s.t.} \qquad \mathbb{E}\{\|w\|_2^2\} \leq \vartheta_2 K. \tag{2.5b}$$

**Theorem 2.2.3.** *The unique solution to* $\mathbf{P}_2$ *in* (2.5) *is given by* $\gamma_2^*(w) = (2\pi\vartheta_2)^{-K/2} \exp(-w^\top w/(2\vartheta_2))$.

Theorem 2.2.3 states that $w(k)$, for all $k \in \mathbb{K}$, are i.i.d. Gaussian random variables with zero mean and co-variance of $\vartheta_2 I$. Upon redefining $\mathfrak{U}(\gamma_1^*) = 1/\mathbb{E}\{\|w\|_2\}$, it can be proved that $\mathfrak{P}(\gamma_2^*)\mathfrak{U}(\gamma_2^*) = 1$. As before, this illustrates the conflict between privacy and utility. Here, it also can be shown that $\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_\infty\} \geq= \sqrt{\vartheta_2}$.

### 2.2.3　Bounded Variations

Both (2.4) and (2.5) allow for somewhat large corruptions to the measurement, albeit with relatively small probabilities that are a function $\vartheta_1$ and $\vartheta_2$. The following problem formulation limits the deviations to a preferred box of size $\varepsilon$:

$$\mathbf{P}_3 : \max_{\gamma \in \Delta(\mathbb{R}^K) \cap \mathcal{C}^1(\mathbb{R}^K, \mathbb{R})} \quad H(w), \tag{2.6a}$$

$$\text{s.t.} \qquad \mathbb{P}\{\|w\|_\infty \leq \varepsilon\} = 1. \tag{2.6b}$$

**Theorem 2.2.4.** *The unique solution to* $\mathbf{P}_3$ *in* (2.6) *is given by* $\gamma_3^*(w) = (2\varepsilon)^{-K}\mathbb{1}_{\|w\|_\infty \leq \varepsilon}$.

Theorem 2.2.4 states that $w(k)$, for all $k \in \mathbb{K}$, are i.i.d. uniform random variables over $[-\varepsilon, +\varepsilon]$. In this case, we can define the utility as $\mathfrak{U}(\gamma_3^*) = 1/(K \max \|w\|_\infty) = 1/(K\varepsilon)$. Note that the utility is a function of $K$ as $K$ entries are perturbed by independent noises and thus the amount of the total perturbation is a function of $K$. We get $\mathfrak{P}(\gamma_3^*)\mathfrak{U}(\gamma_3^*) = 2\varepsilon/(\pi e) = \text{const}$. This again illustrates the conflict between privacy and utility. Here, it also can be shown that $\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_\infty\} \geq \sqrt{2/(\pi e)}\varepsilon$. Thus, it can be ensured that $\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_\infty\} \geq \sqrt{\sigma}$ if $\varepsilon = \sqrt{2/(\pi e \sigma)}$.

### 2.2.4　Numerical Analysis

In this subsection, we investigate the utility-privacy trade-off for the privacy-preserving policies developed using the information-theoretic privacy framework. We may recall that $\mathbb{E}\{\|x_j - \hat{x}_j(z, x_{-j})\|_2^2\} \geq K/(2\pi e) \exp(2H(w)/K)$; see Theorem 2.2.1. Therefore, we use $K/(2\pi e) \exp(2H(w)/K)$ as a proxy for the privacy guarantee. Figure 2.7 (left) illustrates the privacy guarantee versus the expected absolute value $\mathbb{E}\{\|w\|_1\}/K$ for the additive noise in Theorems 2.2.2–2.2.4. As proved in Theorem 2.2.2, the Laplace noise provides the highest privacy guarantee for a fixed level of quality captured by the expected absolute value. However,
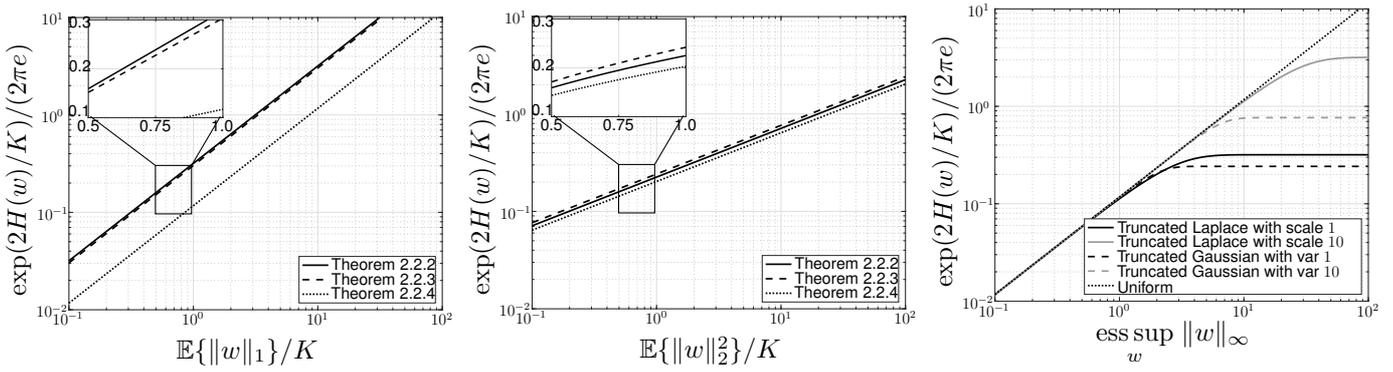
**Figure 2.7:** The privacy guarantee $K/(2\pi e)\exp(2H(w)/K)$ versus (left) the expected absolute value $\mathbb{E}\{\|w\|_1\}/K$ for the additive noise in Theorems 2.2.2–2.2.4, (middle) the variance $\mathbb{E}\{\|w\|_2^2\}/K$ for the additive noise in Theorems 2.2.2–2.2.4, and (right) the essential extrema $\operatorname{ess\,sup}_w\|w\|_\infty$ for the additive noise in Theorems 2.2.2–2.2.4 (with a truncated version being used if the noise has an infinite support).

the same is not true if the measure of quality is inversely proportional with the variance. Figure 2.7 (middle) illustrates the privacy guarantee versus the variance $\mathbb{E}\{\|w\|_2^2\}/K$ for the additive noise in Theorems 2.2.2–2.2.4. Clearly, the Gaussian noise outperforms the rest in terms of the privacy guarantee for a fixed level of quality; see Theorem 2.2.3. Note that in both Figures 2.7 (left) and (middle), Laplace and Gaussian noises perform fairly similarly in terms of privacy guarantees. Finally, neither Gaussian nor Laplace additive noise qualify for the case where the quality is inversely proportional with $\operatorname{ess\,sup}_w\|w\|_\infty$. For a comparison with the uniform noise, in such cases, we should use truncated Gaussian and Laplace noises. Figure 2.7 (right) illustrates the privacy guarantee versus the essential extrema $\operatorname{ess\,sup}_w\|w\|_\infty$ for the additive noise in Theorems 2.2.2–2.2.4 (with a truncated version being used if the noise has an infinite support). As proved in Theorem 2.2.4, the uniform noise is the most privacy-preserving option. Note that, as the size of the additive noise $\operatorname{ess\,sup}_w\|w\|_\infty$ reduces, all continuous density functions approach the uniform density function and thus their privacy guarantees becomes similar.

# 3     Privacy-Preserving Clustering

In this chapter, we consider the problem of clustering in a privacy-preserving manner. First, in Section 3.1, we investigate whether the clusters are privacy-sensitive or not. Subsequently, in Section 3.2, we use locally differentially private versions of smart meter data to extract privacy-preserving clusters.

## 3.1     Are Clusters Privacy Preserving?

Again, consider a private dataset of the form $(x_i(k))_{k \in \mathbb{K}, i \in \mathbb{I}}$ with $\mathbb{K} := \{0, \ldots, K-1\}$ and $\mathbb{I} := \{1, \ldots, n\}$ in which $K \in \mathbb{N}$ (where $\mathbb{N}$ is the set of natural numbers) denotes the time horizon and $n \in \mathbb{N}$ denotes the number of the private entries in the dataset. For instance, $x_i(k)$ may denote the consumption of house $i \in \mathbb{I}$ at time instant $k \in \mathbb{K}$ measured by a smart meter. It is desired to find clusters of $(\mathfrak{C}_\ell)_{\ell=1}^p$, for some positive number $p$, such that $\mathfrak{C}_{\ell_1} \cap \mathfrak{C}_{\ell_2} = \emptyset$ and $\cup_{\ell=1}^p \mathfrak{C}_\ell = \mathbb{I}$. Each cluster represents consumers that have similar consumption patterns. We use the clustering algorithm in [Motlagh et al., 2019]. This method relies on a technique using the delay coordinate embedding for the load-clustering. This is a model-based approach which can handle unequal, asynchronous, and noisy time series. We use the open dataset of smart meter measurements [Department of the Environment and Energy, 2015], also used in [Motlagh et al., 2019] for evaluation of the clustering algorithms and their privacy risks. We particularly use the data of 10,905 homes with at least one full year of time series data. This dataset is referred to as the Smart-Grid Smart-City (SGSC) dataset.

Figure 3.1 illustrates the profile, i.e., empirical probability density function and percentile statistics, of 12 clusters for household energy consumptions in SGSC dataset. Note that the percentile statistics can be computed using the methodology in Chapter 2 to ensure privacy of individuals within the clusters; however, that is out of the scope of this chapter.

In order to understand privacy risks of clustering, we investigate the effort removing an individual household from the SGSC dataset and computing the clusters again. Let $(\mathfrak{C}_\ell^{-j})_{\ell=1}^p$ denote the clustering done on the private dataset $(x_i(k))_{k \in \mathbb{K}, i \in \mathbb{I} \setminus \{j\}}$, i.e., the dataset without the data of household $j$. We expect, if the clusters are privacy preserving, the new clusters do not change in comparison to the old clusters $(\mathfrak{C}_\ell)_{\ell=1}^p$, extracted with the data of household $j$. The total number of changes is given by $(\sum_{\ell=1}^p ||\mathfrak{C}_\ell| - |\mathfrak{C}_\ell^{-j}||)/2 - 1$. Figure 3.2 shows the histogram of the number of the changes in the clusters when removing an individual household from the dataset. Clearly, up to 2.5% of the households can change clusters due to removal of one individual household. This can be caused because of two reasons. First, the clusters might not be privacy preserving, that is, they are highly dependent on the times series of "some" individuals within the dataset. The second reason is that the underlying clustering algorithm is not robust to changes or contains stochastic elements. This might be caused by that the clustering problem is computationally expensive and thus the algorithm can only manage to
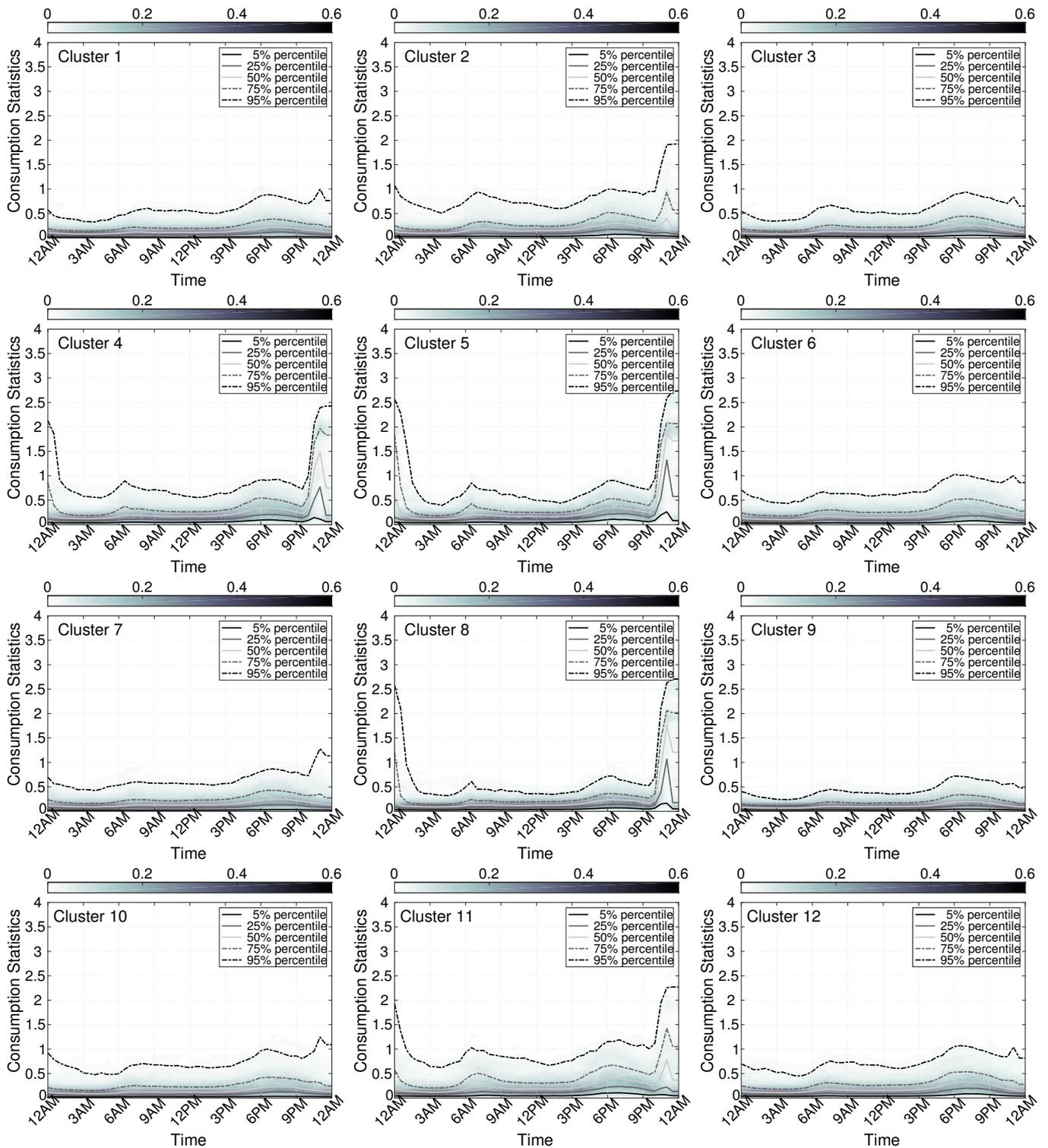
**Figure 3.1: Profile of 12 clusters for household electricity consumptions without noise.**

find locally optimal clusters and, with small changes, the algorithm recovers another locally optimal solution. The algorithm might also contain randomizations to capture better solutions by getting out of suboptimal local solutions. To test this, we can inspect the changes in the clusters when we run the algorithm multiple times without removing households. Figure 3.3 illustrates the histogram of the number of the changes in the clusters without removing households from the dataset with variations caused due to stochasticity in the clustering algorithm. Evidently, there are changes in the clusters; however, the changes, in average, are
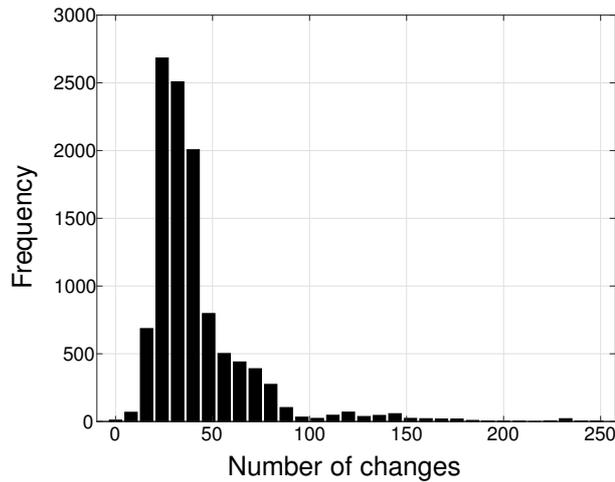
**Figure 3.2: Histogram of the number of the changes in the clusters when removing an individual household from the dataset.**
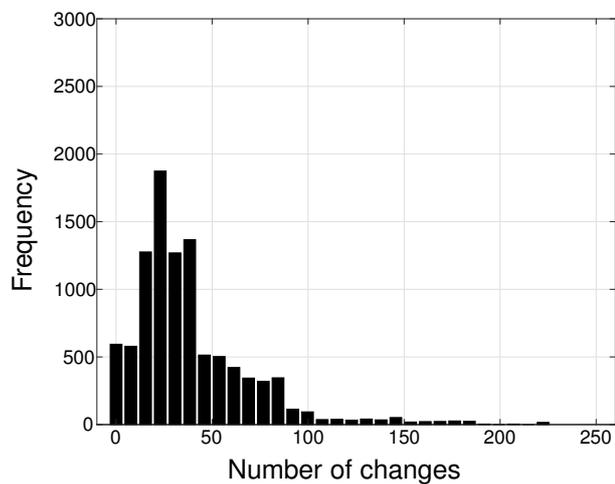


**Figure 3.3: Histogram of the number of the changes in the clusters without removing households from the dataset, variations caused due to stochasticity in the clustering algorithm.**

slightly less than the case in which we remove individual houses from the dataset. Note the higher concentration of the histogram around the smaller changes in Figure 3.3 in comparison to Figure 3.2.

## 3.2    Clustering with Local Differential Privacy

One method to extract a privacy-preserving clustering is to use local differential privacy. To do so, we add noise to the raw data first as

$$y_i(k) = x_i(k) + v_i(k), \qquad \forall k \in \mathbb{K}, \qquad (3.1)$$

where $v_i := (v_i(k))_{k \in \mathbb{K}}$ is a random variable. Since in [Motlagh et al., 2019], the clustering is performed on a delay coordinate embedding of the time series, instead of adding noise to the original load consumptions as above, we add noise the to the parameters of the delay

**Algorithm 1** Optimal correspondence between the clusters from $(\overline{\mathfrak{C}}_\ell)_{\ell=1}^p$ and $(\mathfrak{C}_\ell)_{\ell=1}^p$.

**Require:** $(J(\overline{\mathfrak{C}}_i, \mathfrak{C}_j))_{i,j=1}^p$

**Ensure:** $\sigma$

  $\mathcal{I}_1 \leftarrow \{1, \ldots, p\}$
  $\mathcal{I}_2 \leftarrow \{1, \ldots, p\}$
  **while** $\mathcal{I}_1 \neq \emptyset \wedge \mathcal{I}_2 \neq \emptyset$ **do**
    $(i, j) \leftarrow \arg\max_{(i,j)\in\mathcal{I}_1\times\mathcal{I}_2} J(\overline{\mathfrak{C}}_i, \mathfrak{C}_j)$
    $\sigma(j) \leftarrow i$
    $\mathcal{I}_1 \leftarrow \mathcal{I}_1 \setminus \{i\}$
    $\mathcal{I}_2 \leftarrow \mathcal{I}_2 \setminus \{j\}$
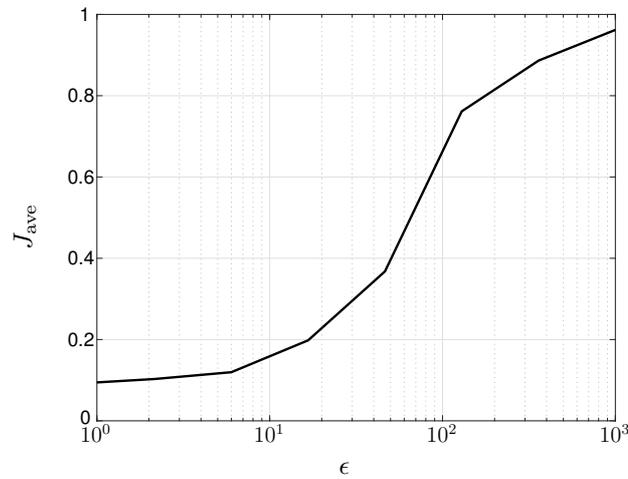  **end while**



**Figure 3.4: Average similarity between clusters with and without privacy-preserving noise versus the differential privacy parameter.**

coordinate embedding of the time series. For each $i \in \mathbb{I}$, $a_i \in \mathbb{R}^m$ is the vector of the parameters of the delay coordinate embedding of the time series $(x_i(k))_{k\in\mathbb{K}}$. We use the additive noise to obfuscate the parameters as

$$\bar{a}_i = a_i + \bar{v}_i,$$

where $\bar{v}_i$ is a vector of i.i.d. Laplace random variables with zero mean and scale $\Delta a/\epsilon$ with $\epsilon$ denoting the differential privacy parameter and $\Delta a$ denoting a constant such that $a_i \in [-\Delta a, \Delta a]^m$ for all $i \in \mathbb{I}$. Let $(\overline{\mathfrak{C}}_\ell)_{\ell=1}^p$ denote the clustering based on noisy parameter vectors $(\bar{a}_i)_{i\in\mathbb{I}}$. We use Jaccard metric to match the corresponding clusters from $(\overline{\mathfrak{C}}_\ell)_{\ell=1}^p$ to $(\mathfrak{C}_\ell)_{\ell=1}^p$. For any pair of clusters $\overline{\mathfrak{C}}_{\ell_1}$ to $\mathfrak{C}_{\ell_2}$, we define the Jaccard metric as

$$J(\overline{\mathfrak{C}}_{\ell_1}, \mathfrak{C}_{\ell_2}) := \frac{|\overline{\mathfrak{C}}_{\ell_1} \cap \mathfrak{C}_{\ell_2}|}{|\overline{\mathfrak{C}}_{\ell_1} \cup \mathfrak{C}_{\ell_2}|}.$$

Evidently, $0 \leq J(\overline{\mathfrak{C}}_{\ell_1}, \mathfrak{C}_{\ell_2}) \leq 1$ and $J(\overline{\mathfrak{C}}_{\ell_1}, \mathfrak{C}_{\ell_2}) = 1$ if and only if $\overline{\mathfrak{C}}_{\ell_1} = \mathfrak{C}_{\ell_2}$. We follow the Algorithm 1 to find the optimal correspondence between the clusters from $(\overline{\mathfrak{C}}_\ell)_{\ell=1}^p$ and $(\mathfrak{C}_\ell)_{\ell=1}^p$. The cluster most similar to $\mathfrak{C}_\ell$ for $1 \leq \ell \leq p$ is given by $\sigma(\ell)$. This allows us to compute an

average Jaccard similarity between the clusters $\overline{\mathfrak{C}}_{\ell_1}$ and $\mathfrak{C}_{\ell_2}$ as

$$J_{\text{ave}} := \frac{1}{p} \sum_{\ell=1}^{p} J(\overline{\mathfrak{C}}_{\sigma(\ell)}, \mathfrak{C}_{\ell}).$$

Figure 3.4 shows the average similarity between clusters with and without privacy-preserving noise, captured by the average Jaccard metric, versus the differential privacy parameter $\epsilon$. As $\epsilon$ increases, i.e., the privacy guarantee weakens, the clusters become more similar.

# 4    Conclusions

Sharing the massive collection of smart meter energy data is highly beneficial for improving commercial applications, government services and researches. However the growing privacy concerns and confidentiality issues impede data custodians sharing the personal and sensitive information available in the energy data collected. Several methods have been developed to overcome the privacy issues by sharing perturbed versions of the data at the cost of utility loss. However, these methods have considered privacy-preserving sharing of the data in an aggregate level, and so far no work has been conducted on releasing percentile statistics, e.g., median of energy data (which is generally time-series data). In addition, the privacy risk of releasing clusters of household energy consumption data has not been investigated in previous studies. In this report, we studied the research problem from two directions.

In the first part, we tackled the privacy-preserving release of percentile statistics of time-series data using global, local, and stochastic differential privacy as well as information-theoretic privacy. In the second part we examined the privacy risks of clustering household energy consumption data and the effect of using local differential privacy for privacy-preserving clustering. We conducted experiments using publicly available real datasets. The results indicate that the local differential privacy outperforms other methods for privacy-preserving releasing of percentile statistics in terms of the privacy-utility trade-off. Moreover, our results show that releasing clusters of household consumption data can compromise privacy and the use of local differential privacy improves privacy-preservation at the cost of significant utility loss. Therefore, further research is required in the area of privacy-preserving clustering of energy consumption data.

# Bibliography

T. Cover and J. Thomas. *Elements of Information Theory*. John Wiley & Sons, 2 edition, 2012. 20

C. Cuijpers and B.-J. Koops. Smart metering and privacy in Europe: lessons from the Dutch case. In *European data protection: coming of age*, pages 269–293. Springer, 2013. 5

Department of the Environment and Energy. Smart-grid smart-city customer trial data: Electricity use interval reading, 2015. Created: 09/09/2015, Updated: 14/08/2018, Accessed: 17/10/2018, https://search.data.gov.au/dataset/ ds-dga-4e21dea3-9b87-4610-94c7-15a8a77907ef/details?q=sgsc. 10, 24

R. Dewri. Local differential perturbations: Location privacy under approximate knowledge attackers. *IEEE Transactions on Mobile Computing*, 12(12):2360–2372, 2013. 6

F. du Pin Calmon and N. Fawaz. Privacy against statistical inference. In *Communication, Control, and Computing (Allerton), 2012 50th Annual Allerton Conference on*, pages 1401–1408. IEEE, 2012. 6

J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Local privacy and statistical minimax rates. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 429–438, 2013. 6

C. Dwork. Differential privacy: A survey of results. In *International Conference on Theory and Applications of Models of Computation*, pages 1–19. Springer, 2008. 5

C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3–4):211–407, 2014. 5, 7, 12, 17

C. Dwork, M. Naor, T. Pitassi, and G. N. Rothblum. Differential privacy under continual observation. In *Proceedings of the 42nd ACM Symposium on Theory of Computing*, pages 715–724. ACM, 2010. 5

G. Eibl and D. Engel. Differential privacy for real smart metering data. *Computer Science - Research and Development*, 32(1):173–182, Mar 2017. 5

R. Fano. *Transmission of Information: A Statistical Theory of Communications*. M.I.T. Press, 1961. 9

F. Farokhi and H. Sandberg. Fisher information as a measure of privacy: Preserving privacy of households with smart meters using batteries. *IEEE Transactions on Smart Grid*, 2017. 6

F. Farokhi, H. Sandberg, I. Shames, and M. Cantoni. Quadratic Gaussian privacy games. In *Decision and Control (CDC), 2015 IEEE 54th Annual Conference on*, pages 4505–4510. IEEE, 2015a. 6

F. Farokhi, I. Shames, and M. Cantoni. Promoting truthful behavior in participatory-sensing mechanisms. *IEEE Signal Processing Letters*, 22(10):1538–1542, 2015b. 5

U. Greveler, B. Justus, and D. Loehr. Multimedia content identification through smart meter power usage profiles. *Computers, Privacy and Data Protection*, 1:10, 2012. 2, 5

R. Hall, A. Rinaldo, and L. Wasserman. Random differential privacy. *Journal of Privacy and Confidentiality*, 4(2):43–59, 2012. 6, 9, 14

P. Kairouz, S. Oh, and P. Viswanath. Extremal mechanisms for local differential privacy. In *Advances in Neural Information Processing Systems*, pages 2879–2887, 2014. 6

K. Kalantari, L. Sankar, and O. Kosut. On information-theoretic privacy with general distortion cost functions. In *Information Theory (ISIT), 2017 IEEE International Symposium on*, pages 2865–2869. IEEE, 2017. 6

Z. Li, T. J. Oechtering, and D. Gündüz. Smart meter privacy based on adversarial hypothesis testing. In *Information Theory (ISIT), 2017 IEEE International Symposium on*, pages 774–778. IEEE, 2017. 6

A. Machanavajjhala, D. Kifer, J. Abowd, J. Gehrke, and L. Vilhuber. Privacy: Theory meets practice on the map. In *Proceedings of the 2008 IEEE 24th International Conference on Data Engineering*, pages 277–286, 2008. 6, 9, 14

P. McDaniel and S. McLaughlin. Security and privacy challenges in the smart grid. *IEEE Security & Privacy*, 3:75–77, 2009. 2, 5

A. Molina-Markham, P. Shenoy, K. Fu, E. Cecchet, and D. Irwin. Private memoirs of a smart meter. In *Proceedings of the 2nd ACM Workshop on Embedded Sensing Systems for Energy-Efficiency in Building*, pages 61–66. ACM, 2010. 2, 5

O. Motlagh, A. Berry, and L. O'Neil. Clustering of residential electricity customers using load time series. *Applied Energy*, 237:11 – 24, 2019. 10, 24, 26

S. J. Pappu, N. Bhatt, R. Pasumarthy, and A. Rajeswaran. Identifying topology of low voltage (LV) distribution networks based on smart meter data. *IEEE Transactions on Smart Grid*, 2017. 5

K. X. Perez, K. Cetin, M. Baldea, and T. F. Edgar. Development and analysis of residential change-point models from smart meter data. *Energy and Buildings*, 139:351 – 359, 2017. 5

B. I. Rubinstein and F. Aldà. Pain-free random differential privacy with sensitivity sampling. In *International Conference on Machine Learning*, pages 2950–2959, 2017. 6, 9, 14

H. Sandberg, G. Dán, and R. Thobaben. Differentially private state estimation in distribution networks with smart meters. In *Decision and Control (CDC), 2015 IEEE 54th Annual Conference on*, pages 4492–4498. IEEE, 2015. 5

J. Shao. *Mathematical Statistics*. Springer Texts in Statistics. Springer-Verlag New York, 2003. 13

O. Tan, D. Gunduz, and H. V. Poor. Increasing smart meter privacy through energy harvesting and storage devices. *IEEE Journal on Selected Areas in Communications*, 31(7):1331–1341, 2013. 6

T. Tanaka, M. Skoglund, H. Sandberg, and K. H. Johansson. Directed information and privacy loss in cloud-based control. In *American Control Conference (ACC), 2017*, pages 1666–1672. IEEE, 2017. 6

T. Tanaka, P. M. Esfahani, and S. K. Mitter. LQG control with minimum directed information: Semidefinite programming approach. *IEEE Transactions on Automatic Control*, 63(1):37–52, 2018. 6

J. W. Tukey. A survey of sampling from contaminated distributions. *Contributions to probability and statistics*, pages 448–485, 1960. 5

J. W. Tukey, F. Mosteller, and D. C. Hoaglin. *Understanding robust and exploratory data Analysis*. Wiley, 1983. 5

Q. Wang, S. R. Kulkarni, S. Verdú, et al. Universal estimation of information measures for analog sources. *Foundations and Trends in Communications and Information Theory*, 5(3): 265–353, 2009. 12

W. Wang, L. Ying, and J. Zhang. On the relation between identifiability, differential privacy, and mutual-information privacy. *IEEE Transactions on Information Theory*, 62(9):5018–5029, 2016. 6

Y. Wang, Q. Chen, T. Hong, and C. Kang. Review of smart meter data analytics: Applications, methodologies, and challenges. *IEEE Transactions on Smart Grid*, 2018. 5

Z. Wang and G. Zheng. Residential appliances identification and monitoring by a nonintrusive method. *IEEE Transactions on Smart Grid*, 3(1):80–92, 2012. 2, 5

G. Wood and M. Newborough. Dynamic energy-consumption indicators for domestic appliances: environment, behaviour and design. *Energy and Buildings*, 35(8):821–841, 2003. 2, 5

## CONTACT US

**t**   1300 363 400
    +61 3 9345 2176
**e**   enquiries@data61.csiro.au
**w**  www.data61.csiro.au

## AT CSIRO WE SHAPE THE FUTURE

We do this by using science and technology to solve real issues. Our research makes a difference to industry, people and the planet.

## FOR FURTHER INFORMATION

Paul Tyler
Data Privacy Team Leader
**t**   +61 2 9490 5908
**e**   Paul.Tyler@data61.csiro.au
**w**  www.data61.csiro.au

Dali Kaafar
Information Security and Privacy Group Leader
**e**   Dali.Kaafar@data61.csiro.au
**w**  www.data61.csiro.au

Hassan Asghar
Senior Researcher
**e**   Hassan.Asghar@data61.csiro.au
**w**  www.data61.csiro.au

Farhad Farokhi
Researcher
**e**   Farhad.Farokhi@data61.csiro.au
**w**  www.data61.csiro.au

Dinusha Vatsalan
Researcher
**e**   Dinusha.Vatsalan@data61.csiro.au
**w**  www.data61.csiro.au

Sirine M'rabet
Researcher
**e**   Sirine.Mrabet@data61.csiro.au
**w**  www.data61.csiro.au