# Market Value of Differentially-Private Smart Meter Data

Saurab Chhachhi
*Imperial College London*
London, UK
saurab.chhachhi11@imperial.ac.uk

Dr. Fei Teng
*Imperial College London*
London, UK
f.teng@imperial.ac.uk

*Abstract*—This paper proposes a framework to investigate the value of sharing privacy-protected smart meter data between domestic consumers and load serving entities. The framework consists of a discounted differential privacy model to ensure individuals cannot be identified from aggregated data, a ANN-based short-term load forecasting to quantify the impact of data availability and privacy protection on the forecasting error and an optimal procurement problem in day-ahead and balancing markets to assess the market value of the privacy-utility trade-off. The framework demonstrates that when the load profile of a consumer group differs from the system average, which is quantified using the Kullback-Leibler divergence, there is significant value in sharing smart meter data while retaining individual consumer privacy.

*Index Terms*—data markets, differential privacy, load forecasting, smart grid, smart meters

## I. INTRODUCTION

Smart metering for domestic consumers is seen as a key enabler in moving towards a more dynamic, cost-effective, cost-reflective and decarbonised electricity network. It provides benefits for both load serving entities (LSE) and consumers through improved billing accuracy, real-time feedback on consumption and enabling innovative business models which harness demand response and dynamic pricing schemes. Access to granular data, such as half-hourly (HH) consumption of individuals, can be used to discern personal information raising issues around privacy and data usage. An overview of potential privacy concerns and smart meter data misuse is provided in [1]. Survey studies have shown that there are a range of attitudes towards smart meter data privacy. Some consumers are happy to share their consumption data, others are willing to share if details on how such data will be used and importantly how it may benefit the system as well as benefit them personally is provided, and those who are reluctant to share data under any circumstances [2].

Currently meter data are collected and processed on an ad-hoc basis by LSEs and then sent to the settlement body. In the UK, the introduction of smart meters will mean that, for settlement purposes, data will be collected by a centralised data company (DCC), with LSEs no longer being part of the process [3]. Moving towards HH settlement will mean that LSEs will have to forecast HH consumption as opposed to

daily volumes. This raises the question as to whether LSEs should have access to individual consumers HH data and how that access should be provided for forecasting purposes. Ofgem, the UK energy regulator, has discussed the use of privacy-preserving mechanisms but there is a lack of understanding as to the costs and benefits of such measures [3].

Extant literature on smart meter privacy-preserving mechanisms has focused on their effect on data utility (i.e. change in forecasting accuracy) [4]. However the resulting impact on energy procurement costs has, to the best of our knowledge, not been investigated. To this end, we propose a framework to assess the smart meter data privacy cost-benefit trade-off for forecasting purposes, simultaneously addressing the specific privacy concerns discussed above. The framework consists of three alternative settlement and forecasting schemes: one in which HH data sharing is mandatory, one in which HH data are not shared and one where HH is shared but is privacy-preserved using differential privacy (DP). To compare the different schemes a forecasting and procurement strategy for a LSE is developed. It consists of an adaptable short-term load forecasting mechanism and an optimal procurement strategy for the LSE in the day-ahead and balancing markets. This paper makes the following contributions:

- Proposes a framework to explicitly link smart meter data sharing to monetary value incorporating privacy concerns.
- Develops a forecasting and procurement strategy for a price-making LSE engaged in the day-ahead and balancing markets within which the privacy-utility trade-off can be assessed.
- Applies the Kullback-Leibler divergence as a potential indicator of data value within this context.
- Presents a case study using actual smart meter data and historical market prices.

In Section II, we review and outline privacy-preserving mechanisms for smart meter data. Section III outlines the different load settlement schemes. The forecasting and procurement model is detailed in Section IV. Section V presents the results of a numerical case study based on actual domestic smart meter data and market prices. Finally, conclusions are drawn and future research directions discussed in Section VI.

## II. Smart Meter Data Privacy

An LSE is required to forecast its consumer group's aggregated load and then purchase sufficient energy. To produce load forecasts the LSE requires historical consumption data of its consumers, which may vary in levels of aggregation or temporal resolution. As historical consumption data can provide significant amounts of personal information about an individual, the development of privacy-preserving mechanisms for releasing smart meter data has been a growing area of research. Privacy-preserving mechanisms can be categorised into the following [5]:

- Cryptographic methods such as encryption.
- Data manipulation which includes spatial aggregation and sampling, anonymization and differential privacy (DP).
- User demand shaping using batteries.

Smart meter data can be used to discern a wide range of attributes of an individual and it is difficult to specify the range and depth of potential data misuse. When defining privacy, for smart meter data, we must think about what specific information an individual wants to keep private. Although encryption, data aggregation and anonymization provide some increased privacy protection they do not guarantee that an individual cannot be identified or prevent user information leakage [4], [6]. DP offers a mechanism to ensure that an individual ($n$) cannot be specifically identified from within a dataset (e.g. an LSE's consumer group) while potentially preserving a high degree of data utility. This approach can protect an individual's privacy even as data analytics and machine learning techniques evolve and new use cases emerge, since the data cannot be attributed to the individual. DP introduces a mathematical framework to define the likelihood of being identified when making a query (in this case aggregated load) on a dataset. A data reporting mechanism is $\epsilon$-differentially private for $\epsilon > 0$ if for any pair of neighbouring datasets $E(t), E'(t)$ (where the two datasets differ by only one individual) and some aggregated output the following holds [6]:

$$\frac{p(\hat{E}(t)|E(t))}{p(\hat{E}(t)|E'(t))} \leq exp(\epsilon) \qquad (1)$$

To achieve this data are obfuscated by introducing Laplacian noise to each individual load profile with 0 mean and $b$ scaling which is given by:

$$b = \frac{\Delta f_t}{\epsilon} \qquad (2)$$

where $\Delta f_t = \frac{\max(E_{n,t}) - \min(E_{n,t})}{N}$, is the global sensitivity of the output (range of the individual loads in a given period $t$), $\epsilon$ is the privacy budget which indicates the risk of being identified and $N$ is the total number of individuals in the dataset.

Most extant literature on the application of DP and its variants to smart meter data have assumed that the privacy budget is fixed and that queries are independent. A detailed review can be found in [4]. However, smart meter data and the resulting queries are continuously generated and updated.

As a result the privacy loss defined by DP is accumulated across each query, requiring the addition of increasing amount of noise to ensure (1) holds. Over time this degrades data quality and renders new data useless [7]. Techniques have been proposed to overcome this limitation such as selective sampling based on time series dynamics [8] which improve performance but are still sensitive to the number of queries and would not guarantee the specified privacy budget over an infinite time horizon. [7] overcomes this, providing a bounded mechanism by introducing the notion of discounted differential privacy (DDP). It draws upon the concept of discounting from economic theory to propose that data further from the past is less sensitive than current data. The resulting noise scaling can be modelled as a function of the privacy budget (as before, $\epsilon$) and the discount rate (a measure of how much one values past data, $\gamma$):

$$b = \frac{\Delta f_t}{\epsilon(1 - \gamma)}, \gamma \in [0, 1) \qquad (3)$$

If one does not place any value on the privacy of past data then $\gamma = 0$ and the mechanism is equivalent to DP whereas if one places high value on the privacy of past data then $\gamma \to 1$ and the required noise tends to infinity.

## III. Domestic Load Forecasting and Settlement

### A. Non-Half-Hourly Settlement (NHHS)

In the absence of HH smart meter data, electricity settlement is based on system-wide daily load coefficients ($DLC$) which are published ahead of time. DLCs are standardised load profiles which specify the amount of annual consumption a specific consumer group (domestic, SME etc.) consumes in a particular half hour. These are generated based on HH measurement taken from a sample of consumers within each consumer group (for details see [9]). An LSE is only required to forecast daily demand ($E^d$) and is therefore insulated from HH changes while still exposed to HH prices (see Table I).

TABLE I
SETTLEMENT SCHEMES

|  | NHHS | HHS - $DLC^{sys}$ | HHS - $E^{hh}$ | HHS – $DDP(\epsilon, \alpha)$ |
|---|---|---|---|---|
| Forecast Input | $E^d, DLC^{sys}$ | $E^d, DLC^{sys}$ | $E^{hh}$ | $E^{hh} + Lap(\epsilon, \alpha)$ |
| Forecast Parameter | $E^d$ | $E^{hh}$ | $E^{hh}$ | $E^{hh}$ |
| Settlement | $E^d DLC^{sys}$ | $E^{hh}$ | $E^{hh}$ | $E^{hh}$ |
| Risk Exposure | $E^d, \lambda^{bal}$ | $E^{hh}, \lambda^{bal}$ | $E^{hh}, \lambda^{bal}$ | $E^{hh}, \lambda^{bal}$ |

### B. Half-Hourly Settlement (HHS)

To assess what the underlying value of sharing HH data would be we present three alternatives (see Table I): a scheme in which data sharing is mandatory i.e. the LSE has access to all its consumers aggregate unaltered HH data (HHS - $E^{hh}$), a scheme where only aggregate daily data are shared (HHS - $DLC^{sys}$) and a scheme where aggregate HH data are shared but is privacy-protected using DDP (HHS - $DDP(\epsilon, \gamma)$). Under all these schemes settlement is based on actual HH consumption but the data available for forecasting purposes

differs. An overview of the dataflows for each scheme is shown in Fig. 1. The next section details the forecasting and procurement models used.

## IV. MODEL DEFINITION

### A. Short-Term Load Forecasting

Artificial Neural Networks (ANN) have been widely used and perform well for short-term forecasting applications and are able to capture both linear and non-linear dependencies [10], [11]. We use a simple ANN consisting of three layers: input layer, hidden layer, and output layer where the hidden layer has four hidden neurons. The following features are considered:

$$X = [W, WD, SP, E_{t-h}, E_{t-h-1}, E_{t-2h+1}, \\ E_{t-2h}, E_{t-3h}] \quad (4)$$

where $E_{t-*h}$ are the lagged/historical load values and h is the number of periods in the day (48), $W$ is the week in the year, $WD$ is the day of the week, and $SP$ is the settlement period. The model is implemented in Python using the MLPRegressor model in Scikit-Learn.

### B. Load Serving Entity (LSE) Procurement Problem

A LSE needs to procure energy to meets its customer group's load by participating in long-term trading, day-ahead and intra-day markets, and settling any imbalances between purchase volumes and actual consumption in the balancing market. In this paper we focus on day-ahead and balancing markets. The LSE's procurement strategy can be formulated as a two-stage risk-constrained stochastic program similar to [12]. We assume the LSE is a price-making market entity in both the day-ahead and balancing market and account for risk-aversion by including the optimisation of the conditional value-at-risk (CVaR). The formulation is as follows:

$$\min_{d^{da}, d^{bal}} \overbrace{\underbrace{\sum_t \lambda_t^{da} d_t^{da}}_{Day-Ahead} + \underbrace{\sum_s \pi_s \sum_t \lambda_{s,t}^{bal} d_{s,t}^{bal}}_{Balancing}}^{Expected\ Cost(\hat{\Omega})} \\ + \underbrace{\beta \left[ \zeta + \frac{1}{1-\alpha} \sum_s \pi_s \eta_s \right]}_{CVaR} \quad (5)$$

s.t.

$$d_t^{da} + d_{s,t}^{bal} = d_t^{fore} + d_{s,t}^{err}, \forall s, \forall t \quad (5a)$$

$$\sum_t \left[ \lambda_t^{da} d_t^{da} + \sum_t \lambda_{s,t}^{bal} d_{s,t}^{bal} \right] - \zeta \leq \eta_s, \forall s \quad (5b)$$

$$\eta_s \geq 0, \forall s \quad (5c)$$

where $d_t^{da}$ and $d_{s,t}^{bal}$ are the volumes procured by the LSE in the day-ahead and balancing market respectively, $\lambda_t^{da}$ and $\lambda_{s,t}^{bal}$ are the market prices, $\pi_s$ is the scenario probability, $\beta$ is the risk-aversion factor, $\alpha$ is the CVaR confidence interval, $\zeta$ and $\eta_s$ are auxiliary variables to calculate CVaR, $d_t^{fore}$ and $d_{s,t}^{err}$ are
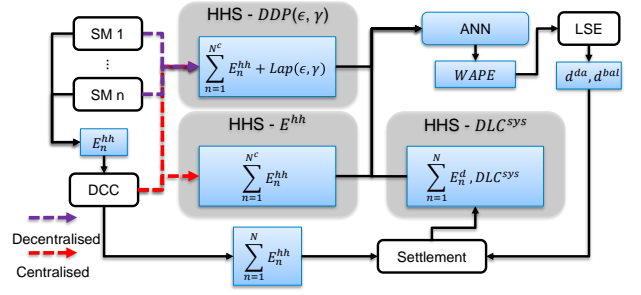


Fig. 1. Overview of Proposed Framework.[1]

the day-ahead forecast load and the realised error. Given that the LSE is a price-making entity $\lambda_t^{da}$ and $\lambda_{s,t}^{bal}$ are dependent on the total demand. These can be modelled as piece-wise linear curves:

$$\lambda_t^{da} = \sum_b \lambda_b^G u_{t,b}^{da}, \forall t \quad (5d)$$

$$\lambda_{s,t}^{bal} = \sum_f \lambda_b^F u_{s,t,f}^{bal}, \forall s, \forall t \quad (5e)$$

$$D_t^{sys} - \frac{\Delta}{2} \leq \sum_b u_{t,b}^{da} \tilde{D}_b^{sys} \leq D_t^{sys} + \frac{\Delta}{2}, \forall t \quad (5f)$$

$$D_{s,t}^{imb} - \frac{\Delta}{2} \leq \sum_f u_{s,t,f}^{bal} \tilde{D}_f^{imb} \leq D_{s,t}^{imb} + \frac{\Delta}{2}, \forall s, \forall t \quad (5g)$$

where $D_t^{sys} = D_t^{da} + d_t^{da}$, the total system demand in period $t$ and $D_{s,t}^{imb} = D_{s,t}^{bal} + d_{s,t}^{bal}$, the total system imbalance in scenario $s$, $\tilde{D}_b^{da}$ is a discretisation of the system demand into $b$ increments of $\Delta$ (similarly for $\tilde{D}_f^{bal}$) and $u_{t,b}^{da}$ and $u_{s,t,f}^{bal}$ are binary variables which select the appropriate demand level. The resulting model is a MIQP due to the products of binary and continuous variables ($u_{t,b}^{da} d_t^{da}$ and $u_{s,t,f}^{bal} d_{s,t}^{bal}$). These can be linearised by replacing the bilinear terms with a new variable on which a number of constraints are imposed giving an exact MILP reformulation [14]. For example the term $u_{t,b}^{da} d_t^{da}$ can replaced by an auxiliary variable, $C_{t,b}^{da}$, and four additional constraints:

$$C_{t,b}^{da} \leq u_{t,b}^{da} d_{t,b}^{max}, \forall t, \forall b \quad (5h)$$

$$C_{t,b}^{da} \leq d_t^{da}, \forall t, \forall b \quad (5i)$$

$$C_{t,b}^{da} \geq d_t^{da} - (1 - u_{t,b}^{da}) d_{t,b}^{max}, \forall t, \forall b \quad (5j)$$

$$C_{t,b}^{da} \geq 0, \forall t, \forall b \quad (5k)$$

The MILP reformulation is implemented in FICO™ Xpress through the Python API.

### C. Assessment Metrics

To gauge the difference between the load profile of the LSE's consumer group and the rest of the system we employ the Kullback-Leibler divergence (KLD). Its value can be interpreted as the information gain achieved if HH data from the consumer group ($DLC^c$) is used instead of system level data

[1]Laplacian noise can be constructed by summing $n$ i.i.d. Gamma distributions allowing for decentralised noise addition at the smart meter [13].

$(DLC^{sys})$. We assume the average weekly $DLC$ variations are normally distributed. In this context it can be defined as follows [15]:

$$KLD = \sum_{t \in T^w} \left[ \log \frac{\sigma_t^{sys}}{\sigma_t^c} + \frac{\sigma_t^{c2} + (\mu_t^{sys} - \mu_t^c)^2}{2\sigma_t^{sys2}} - \frac{1}{2} \right] \quad (6)$$

where $\mu$ and $\sigma$ are the mean and standard deviation of the $DLC$ for each HH of the week respectively.

To measure the accuracy of the forecasts we employ the weighted-absolute percentage error (WAPE) metric as it exhibits more stable behaviour for values close to zero. This is especially relevant in the presence of PV and battery storage as net load profiles can be negative.

$$WAPE = \frac{\sum_t \left| E_t^{act} - E_t^{fore} \right|}{\sum_t E_t^{act}} \quad (7)$$

## V. CASE STUDY

### A. Data

We use smart meter data from the CER Behavioural Trials which includes 6010 residential and SME consumers for a period of 75 weeks [16][2]. For 50% of consumers synthetic PV [17] and EV [18] load profiles are added to better reflect the increased load diversity expected in the future. Day-ahead and balancing market bidding curves are generated based on historical UK market data for 2018 from Elexon. The WAPEs of the ANN described in Section IV-A are used to 50 generate demand forecast scenarios, assuming they are normally distributed and then scaled to a representative UK system level based on the share of total load of consumers in the CER dataset.

Fig. 2a shows the KLDs for randomly sampled meters under various proportions of consumers. When the meters are a small proportion of the total consumers the KLDs of the aggregated load can be large but as the proportion increases the KLDs decrease significantly. We argue that as LSEs begin to offer more innovative and targeted tariff mechanisms KLDs could be large even for large groups of consumers as they change consumption patterns based on particular time-varying incentives. We select four groups based on K-Means clustering of average DLC for each consumer across the week (shown in grey in Fig. 3) to test the framework. To add context we plot the resulting clusters on Fig. 2a (A-D). From Fig. 2b, which shows the forecast error, it is clear that as the KLD of a consumer group increases, having HH data for that specific group results in a greater reduction in forecasting error.

### B. Results

*1) Forecast Accuracy for Different Consumer Group:* The top plots in Fig. 3 show the weekly average DLCs for the selected consumer groups. The bottom plots show the WAPE under each scheme. It is clear that when the average DLC of

---

[2]After data processing to remove periods and meters for which less than 95% of data was available.



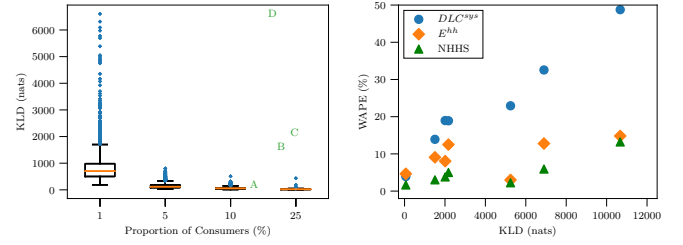(a) KLD for Different Market Shares   (b) WAPE for Select KLD

Fig. 2.  KLD and WAPE for 2000 Sampled Groups.

the consumer group is similar to the system (Group A), reflected in a low KLD, the WAPE is low even in the case where HH data are not provided ($DLC^{sys}$). As a result providing HH data using the DDP mechanism does not increase utility. However, as the group KLD increases, there is an increase in the difference between the WAPE with HH data and the WAPE without access to HH data. This provides a range within which an LSE is able to explore the privacy-utility trade-off. For example for Group C the LSE would be able to gain a 5% reduction in WAPE while providing a $(\epsilon = .25, \gamma = .75)$ level of privacy. This shows that privacy-preservation can be achieved without significantly degrading data utility.

*2) Market Value:* Fig. 4a shows the procurement costs based on scenarios generated for the different settlement schemes for Group C. An LSE can make significant cost savings while still providing consumers with privacy. On average we see that a 1% increase in WAPE results in a 0.8 - 1% cost reduction. A greater reduction is observed in the CVaR as a 1% increase in WAPE results in a 2-3% reduction in CVaR. The value of better forecasting accuracy and hence HH data are also highly dependent on the market dynamics. At peak times uncertainty is more expensive, as there is less flexibility in the balancing market when overall demand is high, resulting in larger cost differences between the schemes.

*3) Heterogeneous Privacy Preferences:* As privacy concerns vary, we investigate how the costs change when only a fraction of the consumer group has privacy concerns. Assuming a proportion of the consumer group $p$ has privacy concerns we generate forecasts separately for them using $DDP$, with the load for the remaining consumers modelled using $E^{hh}$. Fig. 4b shows the resulting procurement costs using this method for $DDP(.25, .75)$. The dots represent the weighted average cost ($\hat{\Omega}^{exp}$) that would be expected based on the proportion $p$. Splitting consumers based on privacy concerns can improve overall data quality and reduce overall procurement costs when $p$ is low. However a trade-off is observed between reduced data degradation, as less noise is added, and benefits of aggregation, which smoothens the load profile.
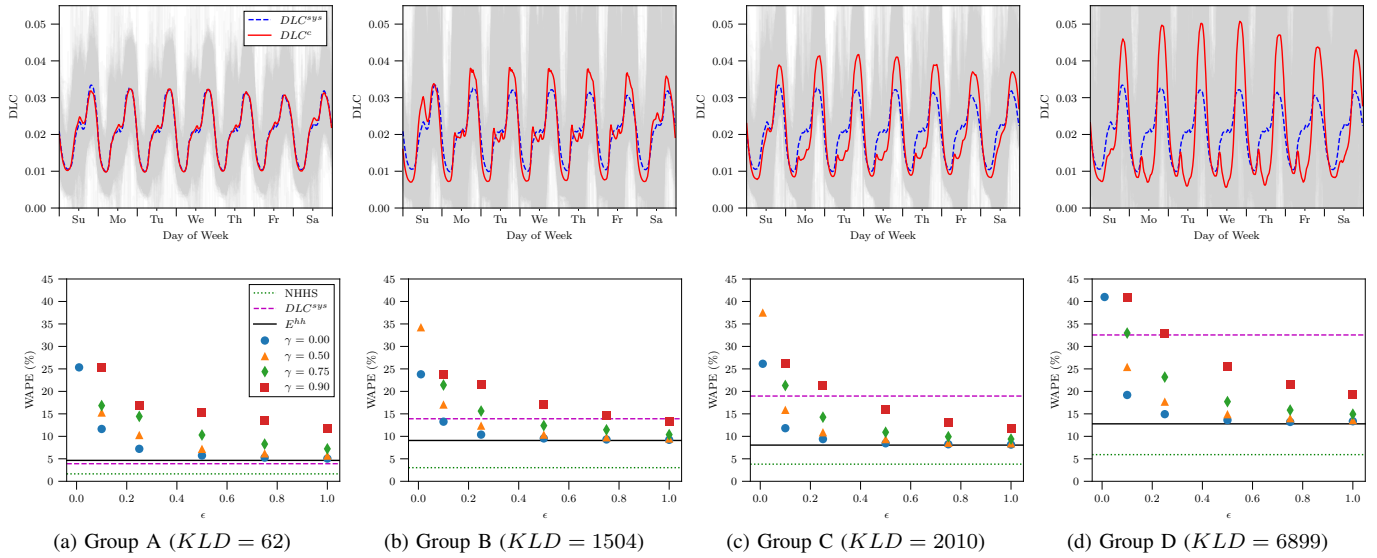
Fig. 3. Group Weekly Average DLC (top), Forecasting Error for Different Settlement Mechanisms (bottom).
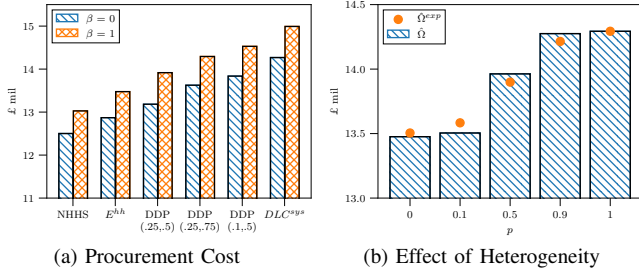


Fig. 4. Procurement Cost for the Different Schemes.

## VI. CONCLUSIONS AND FUTURE WORK

This paper investigated the value of sharing smart meter data applying a framework consisting of a discounted differential privacy model to ensure individuals cannot be identified from aggregated data, a short-term load forecasting method using ANN to quantify the impact of data availability and privacy protection on the forecasting error, and an optimal procurement problem, to assess the market value of the privacy-utility trade-off introduced by DDP. Results show that when the load profile of a LSE's consumer group differs from the system average, which is increasingly relevant with the introduction of dynamic tariffs, and distributed storage and generation, there is significant value in sharing data while retaining individual consumer privacy. Further work is needed to assess how the benefits of smart meter data sharing can be distributed through the development of privacy differentiated tariffs or data markets to incentivise data sharing. In addition, reducing the global sensitivity parameter ($\Delta f$) by optimising the noise introduced by the DDP mechanism and explicitly incorporating heterogeneous privacy preferences would further improve performance. This framework could also be extended to forecasting demand response, net load and flexibility as well as including additional data streams such as user preferences and appliance information.

## REFERENCES

[1] C. Véliz and P. Grunewald, "Protecting data privacy is key to a smart energy future," *Nature Energy*, vol. 3, no. 9, pp. 702–704, 2018.
[2] A. Dickman and A. P. Aslaksen, "Consumer attitudes to DNO access to half hourly electricity consumption data," Ipsos Mori, Tech. Rep., 2017.
[3] Ofgem, "Electricity retail market-wide half-hourly settlement: consultation," Tech. Rep., 2020.
[4] M. U. Hassan, M. H. Rehmani, and J. Chen, "Differential privacy techniques for cyber physical systems: A survey," *IEEE Commun. Surveys Tuts*, vol. 22, no. 1, pp. 746–789, 2020.
[5] G. Giaconi, D. Gunduz, and H. V. Poor, "Privacy-Aware Smart Metering: Progress and Challenges," *IEEE Signal Process. Mag.*, 2018.
[6] C. Dwork and A. Roth, "The Algorithmic Foundations of Differential Privacy," *Found. Trends Theor. Comput. Sci.*, vol. 9, no. 3–4, pp. 211–407, 2014.
[7] F. Farokhi, "Temporally discounted differential privacy for evolving datasets on an infinite horizon," in *Proc. 11th ACM/IEEE Int. Conf. on Cyber-Physical Syst.*, 2020.
[8] J. W. Kim, B. Jang, and H. Yoo, "Privacy-preserving aggregation of personal health data streams," *PLoS ONE*, 2018.
[9] Elexon, "Load Profiles and their use in Electricity Settlement," Tech. Rep., 2018.
[10] Y. Wang, Q. Chen, M. Sun, C. Kang, and Q. Xia, "An Ensemble Forecasting Method for the Aggregated Load with Subprofiles," *IEEE Trans. on Smart Grid*, 2018.
[11] S. I. Vagropoulos, E. G. Kardakos, C. K. Simoglou, A. G. Bakirtzis, and J. P. S. Catalão, "Artificial neural network-based methodology for short-term electric load scenario generation," in *Proc. 18th Int. Conf. on Intell. Syst. Appl. to Power Syst.*, 2015, pp. 1–6.
[12] M. Song and M. Amelin, "Price-Maker Bidding in Day-Ahead Electricity Market for a Retailer With Flexible Demands," *IEEE Trans. on Power Syst.*, vol. 33, no. 2, pp. 1948–1958, 2018.
[13] G. Ács and C. Castelluccia, "I have a DREAM! (DiffeRentially privatE smArt Metering)," in *Lecture Notes in Comput. Sci.*, 2011.
[14] H. P. Williams, *Building integer programming models I.* Wiley, 2013.
[15] S. J. Roberts and W. D. Penny, "Variational bayes for generalized autoregressive models," *IEEE Trans. on Signal Process.*, vol. 50, no. 9, pp. 2245–2257, 2002.
[16] "CER Smart Metering Project - Electricity Customer Behaviour Trial, 2009-2010." Irish Social Science Data Archive, 2012. [Online]. Available: www.ucd.ie/issda/CER-electricity
[17] S. Pfenninger and I. Staffell, "Long-term patterns of European PV output using 30 years of validated hourly reanalysis and satellite data," *Energy*, vol. 114, pp. 1251–1265, 2016.
[18] T. Dodson and S. Slater, "Electric Vehicle Charging Behaviour Study, Final Report for National Grid ESO," Element Energy, Tech. Rep., 2019.