

Data privacy issues in the AMI

ICT Support for Adaptiveness and (Cyber)security in the Smart Grid

DAT300/DIT668

Valentin Tudor
Computer Science and Engineering Department
Chalmers University of Technology
Gothenburg, Sweden

22nd of September 2016



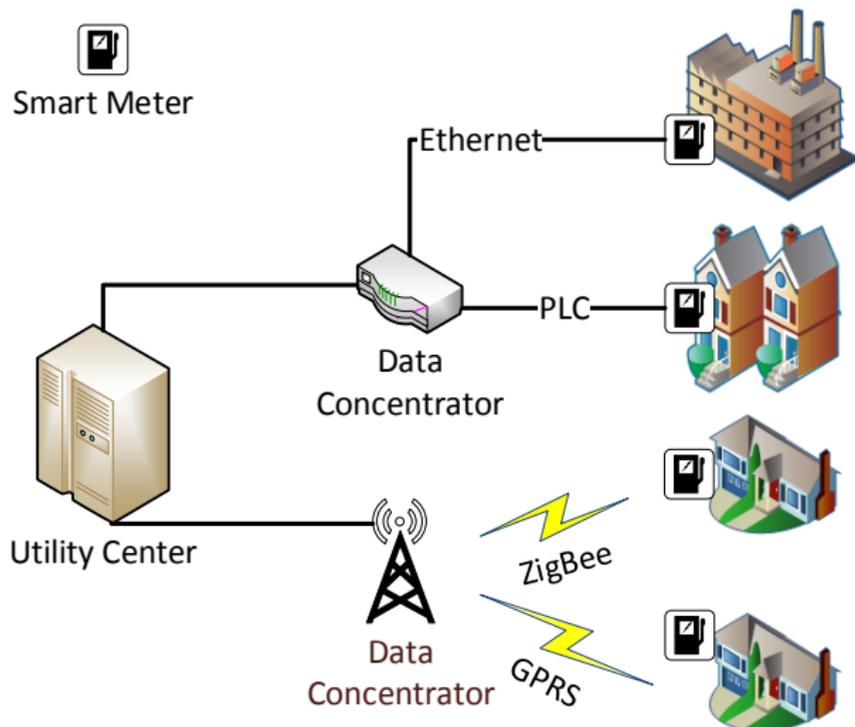
CHALMERS
UNIVERSITY OF TECHNOLOGY



UNIVERSITY OF GOTHENBURG

- The Advanced Metering Infrastructure (AMI)
- AMI data - utility and privacy
- Privacy issues - de-anonymization and de-pseudonymization
- Differential-privacy and AMI data
- AMI data application - energy load forecast using DP-aggregated AMI data
- Conclusion

The Advanced Metering Infrastructure (AMI)



Collecting energy related data



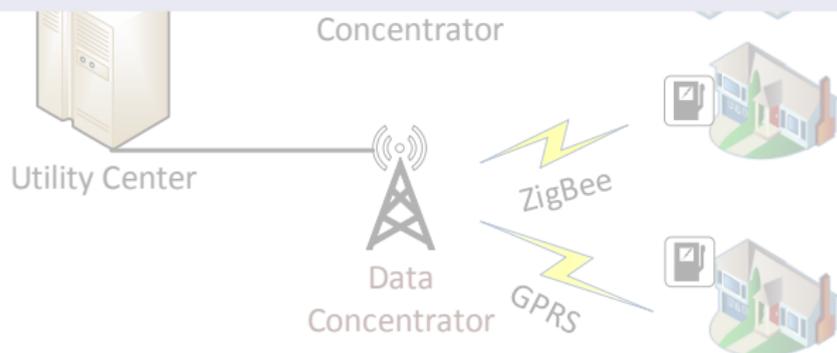
Smart Meter

Ethernet



Before AMI

- data collection frequency - monthly
- only energy consumption readings
- no direct communication to the meter



Collecting data from AMI



Smart Meter

Ethernet



Before AMI

- data collection frequency - monthly
- only energy consumption readings
- no direct communication to the meter



Concentrator



With AMI

- data collection frequency - minutes or less
- energy consumption, power quality info., alarms
- two-way communication with the meter

Utilizing large data from AMI

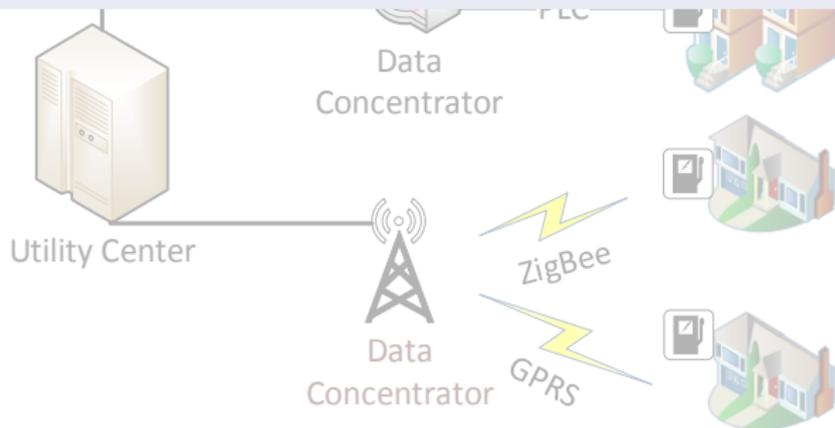


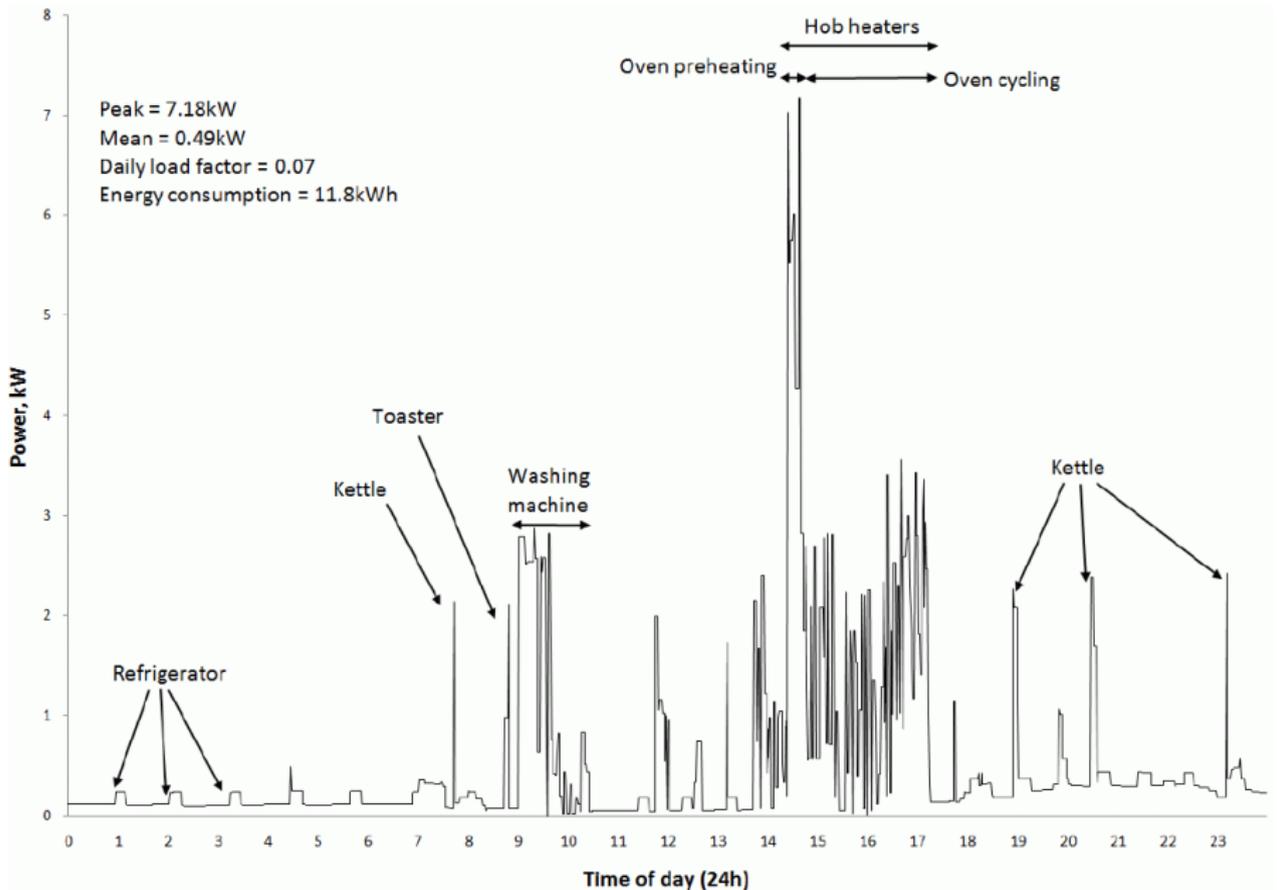
Smart Meter

Ethernet

Many opportunities - great potential

- billing, grid operation, marketing, demand-response, micro-markets, load forecast





Utilizing large data from AMI



Smart Meter

Ethernet



Many opportunities - great potential

- billing, grid operation, marketing, demand-response, micro-markets, load forecast

Privacy challenges

- fine-grained AMI data - sensitive information^a
- privacy invasive techniques against anonymization^{b,c,d}

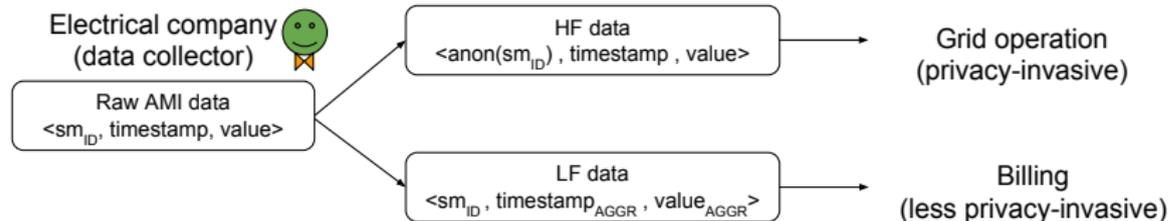
^aMármol, F.G., Sorge, C., Ugus, O., and Pérez, G. M. 2012. Do not snoop my habits: preserving privacy in the smart grid. In IEEE Communications Magazine, 50(5), pp.166-172.

^bJawurek, M., Johns, M. and Rieck, K., 2011, December. Smart metering de-pseudonymization. In Proceedings of the 27th Annual Computer Security Applications Conference (pp. 227-236). ACM.

^cTudor, V., Almgren, M. and Papatriantafidou, M., 2013. Analysis of the impact of data granularity on privacy for the smart grid. In Proceedings of the 12th ACM workshop on Workshop on privacy in the electronic society.

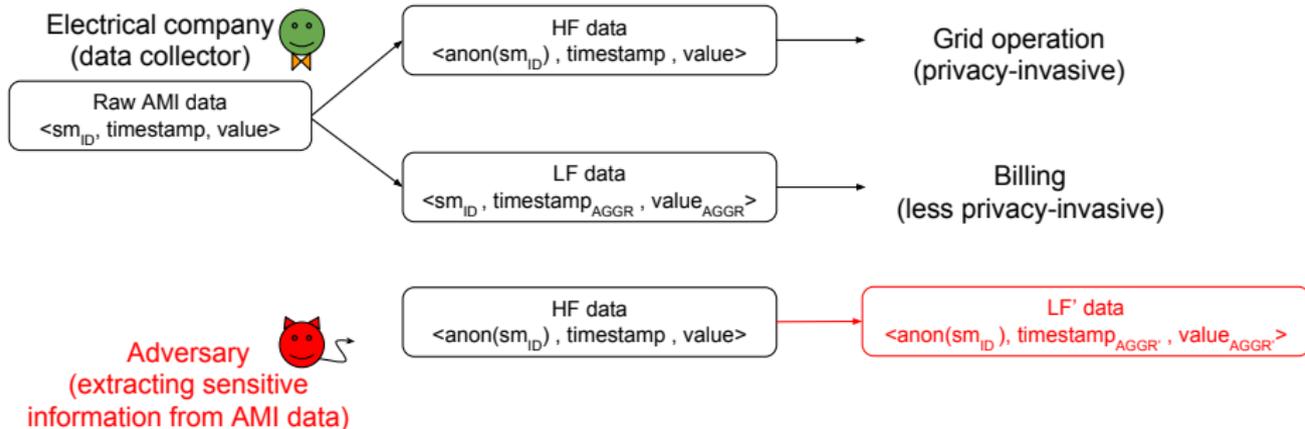
^dTudor, V., Almgren, M., and Papatriantafidou, M., 2015. A study on data de-pseudonymization in the smart grid. In Proceedings of the Eighth European Workshop on System Security

AMI data anonymization - remove PII¹ and aggregation

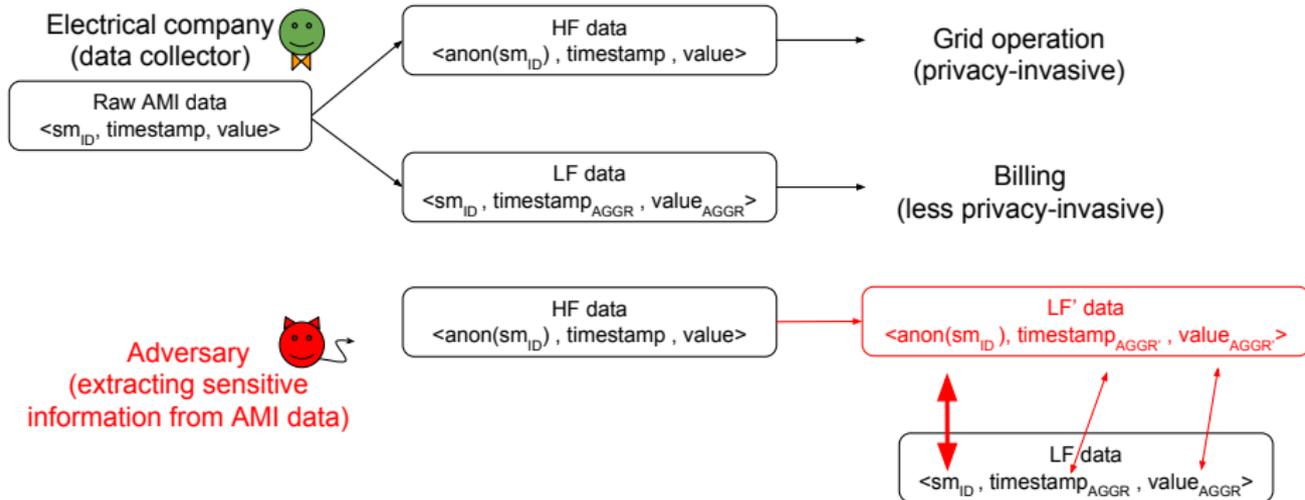


¹ PII - Personal Identifiable Information

De-anonymization



De-anonymization



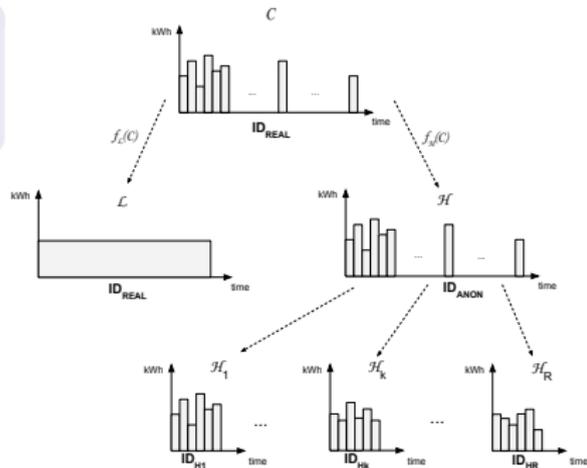
The adversary de-anonymizes AMI datasets using information extracted from AMI data itself.

A bit of formalism

$C = \{(identifier, timestamp, value)\}$ - original dataset

There are two functions such that two new datasets can be derived

$$\begin{cases} \mathcal{H} & = & f_H(C) \\ \mathcal{L} & = & f_L(C) \end{cases}$$



De-anonymization

$C = \{(identifier, timestamp, value)\}$ - original dataset

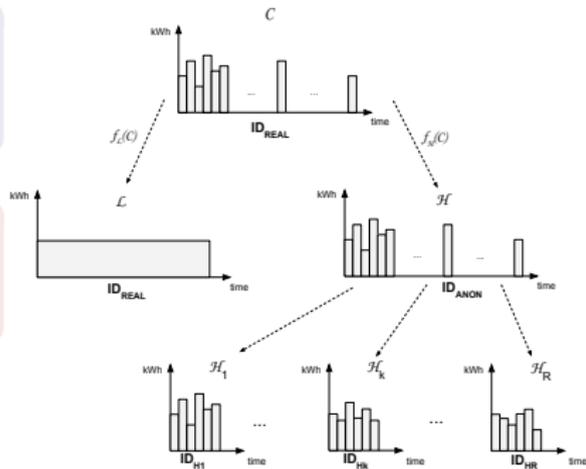
There are two functions such that two new datasets can be derived

$$\begin{cases} \mathcal{H} &= f_H(C) \\ \mathcal{L} &= f_L(C) \end{cases}$$

An adversary is interested in matching the identities in \mathcal{H} with \mathcal{L} .

There is a function $g(\cdot)$, such that $\mathcal{L}' = g(\mathcal{H})$.

Link entries $\mathcal{L}' \sim \mathcal{L}$.



De-pseudonymization

$C = \{(identifier, timestamp, value)\}$ - original dataset

There are two functions such that two new datasets can be derived

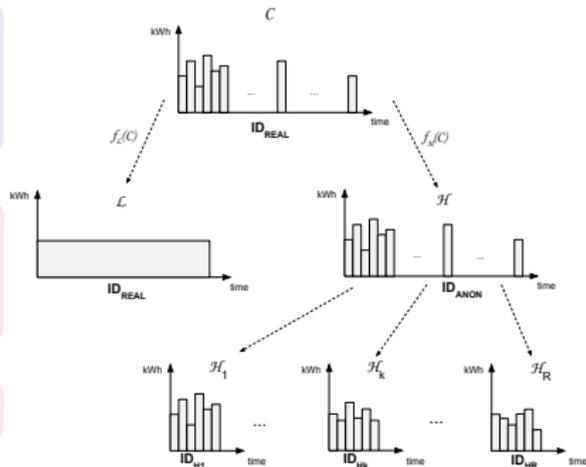
$$\begin{cases} \mathcal{H} &= f_H(C) \\ \mathcal{L} &= f_L(C) \end{cases}$$

An adversary is interested in matching the identities in \mathcal{H} with \mathcal{L} .

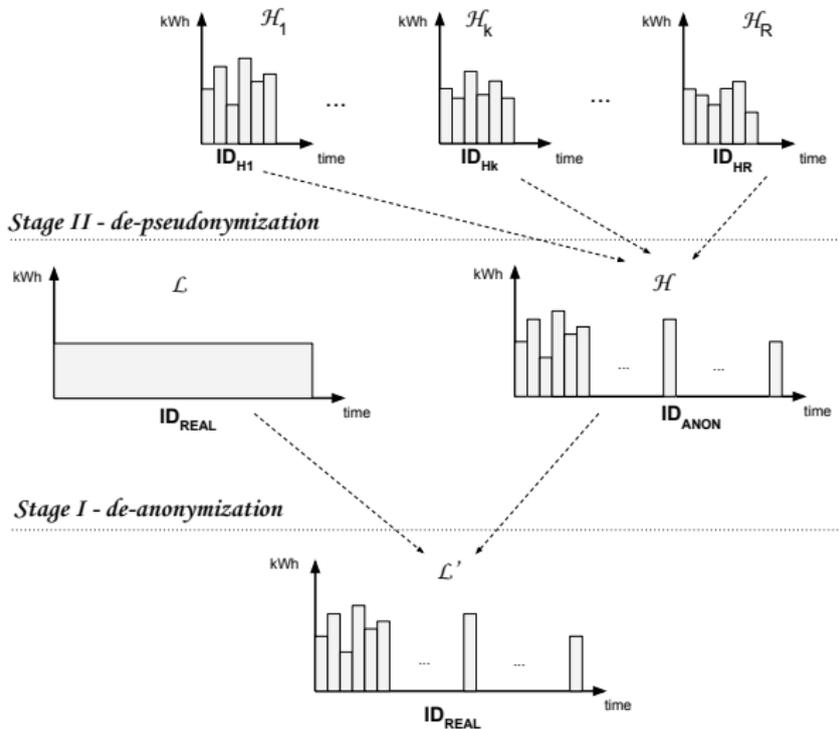
There is a function $g(\cdot)$, such that $\mathcal{L}' = g(\mathcal{H})$.

Link entries $\mathcal{L}' \sim \mathcal{L}$.

There is a function $r(\cdot)$, such that $\mathcal{H}_{k,k+1} = r(\mathcal{H}_k, \mathcal{H}_{k+1})$

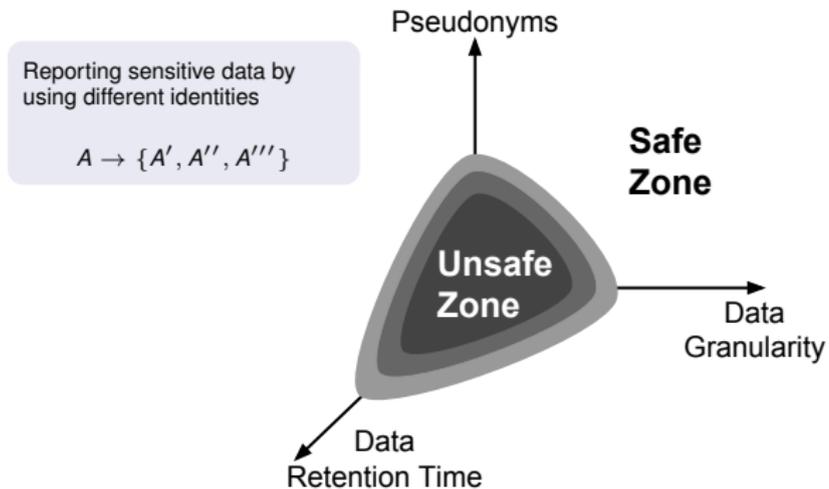


Complete adversarial picture

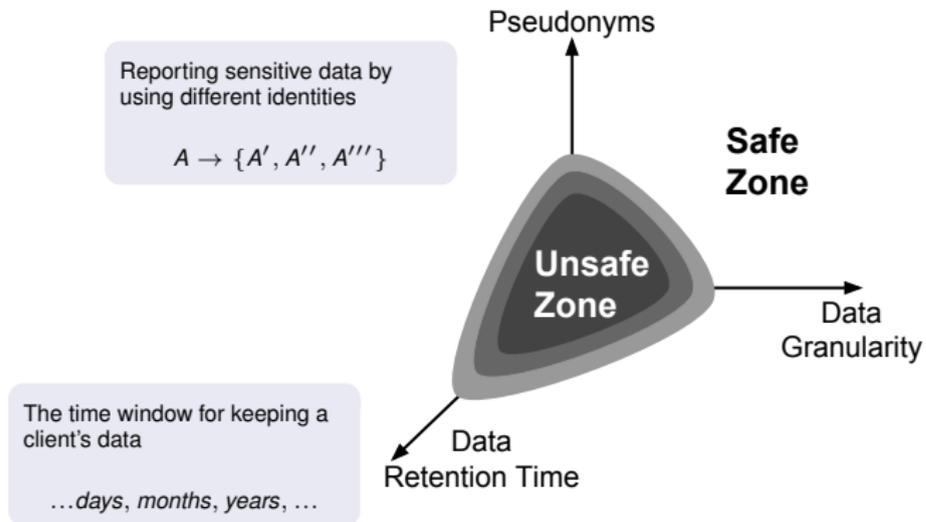


- The Advanced Metering Infrastructure (AMI)
- AMI data - utility and privacy
- **Privacy issues - de-anonymization and de-pseudonymization**
- Differential-privacy and AMI data
- AMI data application - energy load forecast using DP-aggregated AMI data
- Conclusion

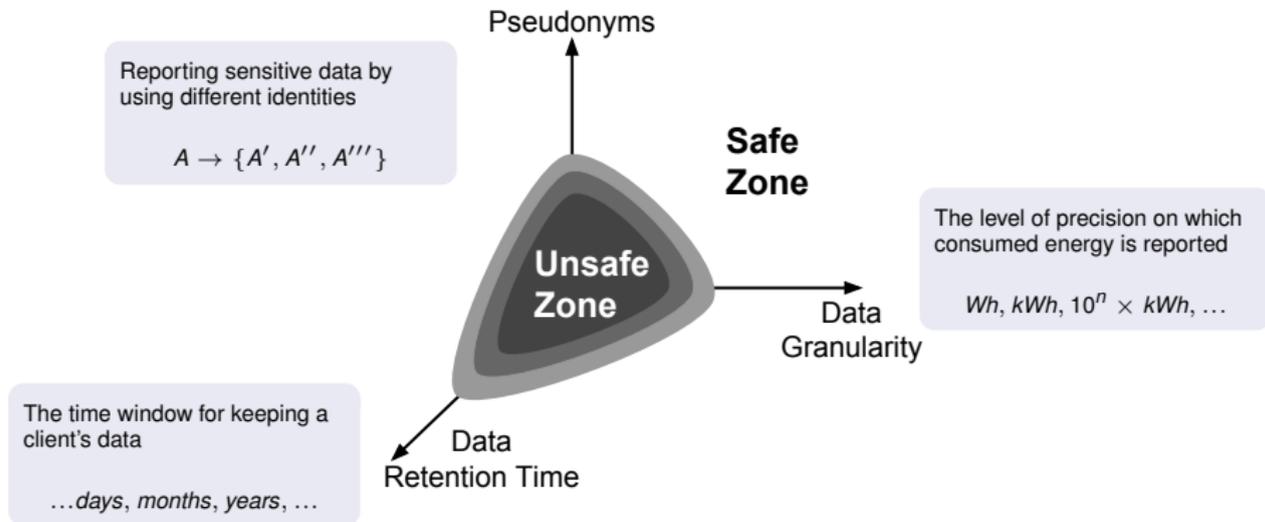
AMI data characteristics



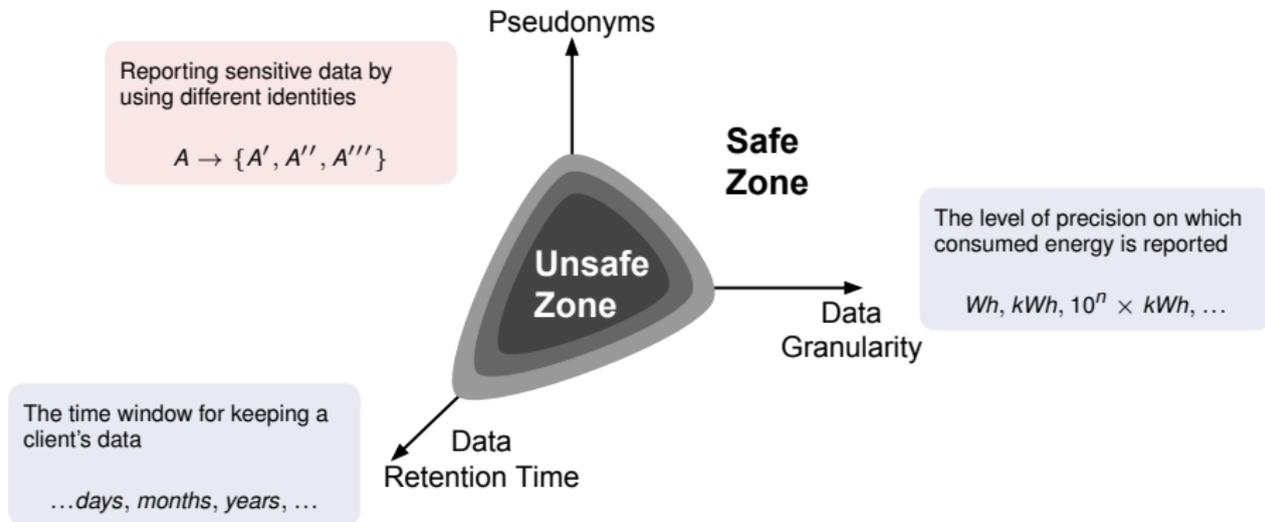
AMI data characteristics



AMI data characteristics



AMI data characteristics



How do the properties of the **HF** datasets influence the effectiveness of utilizing pseudonyms?

Linking together data produced by the same source (household), but stored under different pseudonyms.

We assume a scenario in which an adversary:

- 1 Gets hold of two HF datasets.
- 2 For each household computes a number of features based on the data in the two datasets.
- 3 Uses the features to link together the identities used in the two HF datasets.
- 4 For each correctly linked identity the adversary obtains an extended HF dataset.

De-pseudonymization problem

Linking together data produced by the same source (household), but stored under different pseudonyms.

We assume a scenario in which an adversary:

- 1 Gets hold of two HF datasets.
- 2 For each household computes a number of features based on the data in the two datasets.
- 3 Uses the features to link together the identities used in the two HF datasets.
- 4 For each correctly linked identity the adversary obtains an extended HF dataset.

What is the **effect** of the dataset **size** and collection **season** on the **de-pseudonymization ratio**?

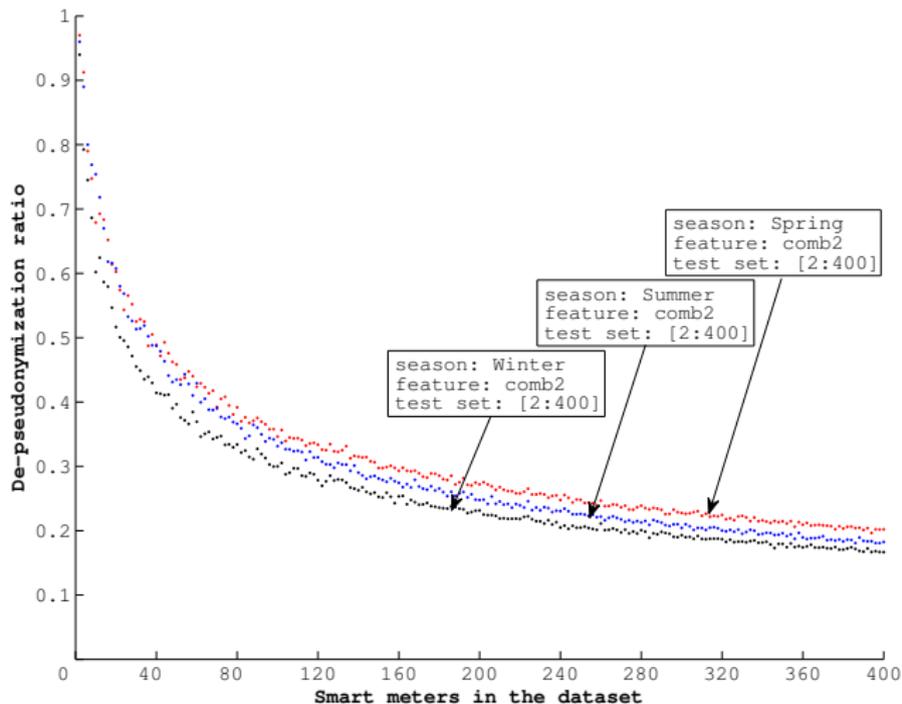
Choosing features

Statistical features that can be computed efficiently on HF data:

#	Feature name	Abbrev.	Description
1	Standard deviation	std	variation of energy consumption
2	Mode	mode	most common consumption value
3	Mean consumption	meanc	average energy consumption
4	Max consumption	maxc	maximum energy consumption
5	Coefficient of variation	cv	ratio of standard deviation to mean

Each household data record becomes a point in a multi-dimensional space → a distance metric can be used to find similar households.

Seasonal results - combined features std, mode, meanc



Random guessing

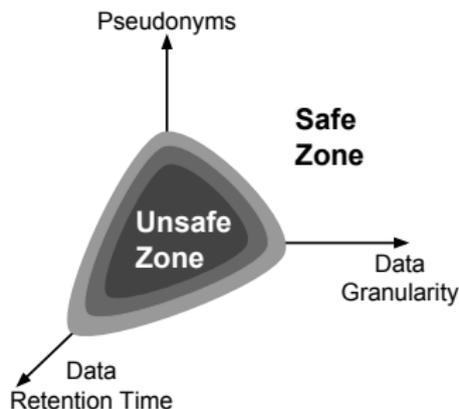
Assume that we have two HF datasets using different pseudonyms for the N households stored in each.

Probability to match:

- 1 pseudonym: $P[1 \text{ out of } N \text{ matched}] = \frac{1}{N}$
- 2 pseudonyms: $P[2 \text{ out of } N \text{ matched}] = \frac{1}{N} \times \frac{1}{N-1}$
- ...
- k pseudonyms:
 $P[k \text{ out of } N \text{ matched}] = \frac{1}{N} \times \frac{1}{N-1} \times \dots \times \frac{1}{N-k+1}$
 $P[k \text{ out of } N \text{ matched}] = \prod_{l=1}^k \frac{1}{N-l+1}$
- N pseudonyms: $P[N \text{ out of } N \text{ matched}] = \frac{1}{N!}$

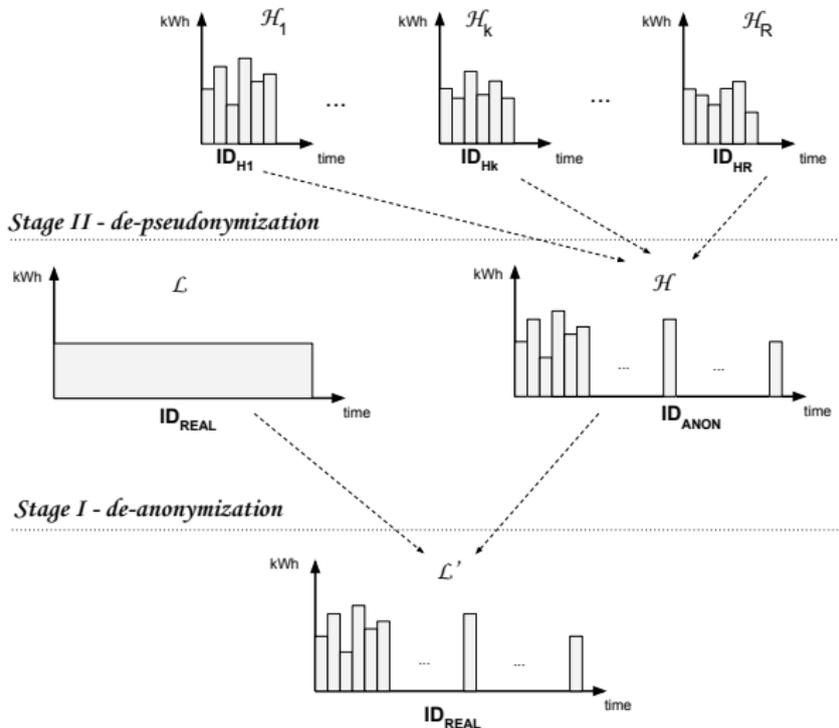
Random guessing might work for datasets with a small number of households, but it becomes harder as the size of datasets increases.

Summary - de-pseudonymization

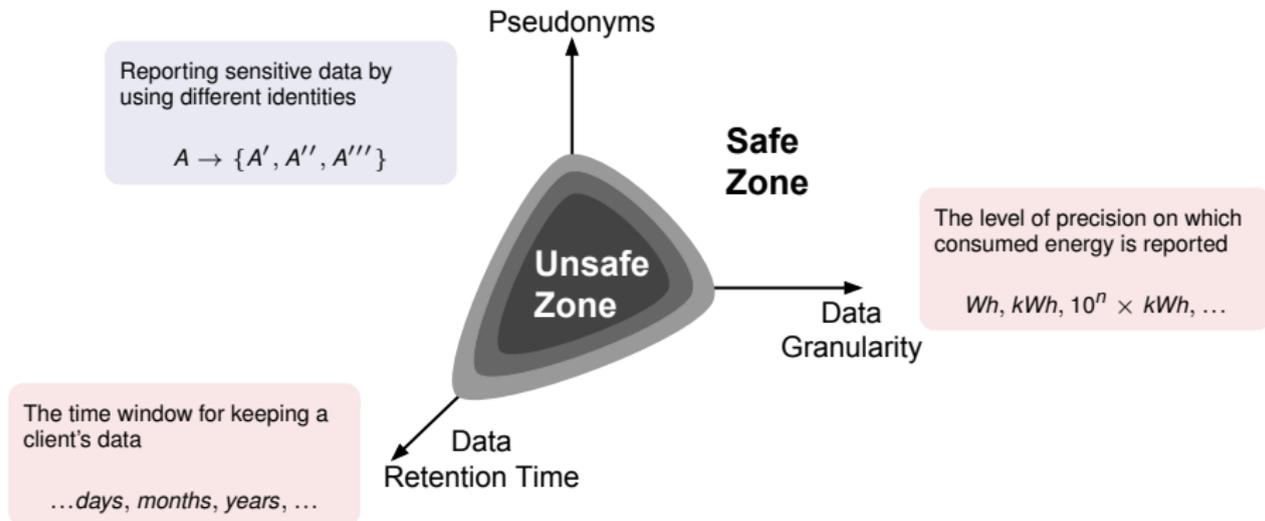


- The **number of households** in the dataset and **collection season** influences the efficiency of the de-pseudonymization.
- The number of re-identified households is not proportional with the size of the dataset.
- The **characteristics** of the Advanced Metering Infrastructure **dataset** should be taken into consideration when evaluating or developing Privacy Enhancing Technologies for this domain.

Complete adversarial picture



AMI data characteristics



Adversarial strategy

1. The adversary gets hold on two datasets, one LF and one HF.
2. For each round (time period of data):
 - the adversary identifies unique smart meters based on their energy consumption index values
 - the values for the identified smart meters are removed from the future rounds
3. The adversary repeats this for each round until she has identified all smart meters or she has used all time periods of data.

Her purpose is to identify uniquely a large number of customers.

$HF \rightarrow LF' \rightarrow LF$

Probabilistic model - Game of balls and bins

Begin

The bins are created based on their width (**data granularity**).



Probabilistic model - Game of balls and bins

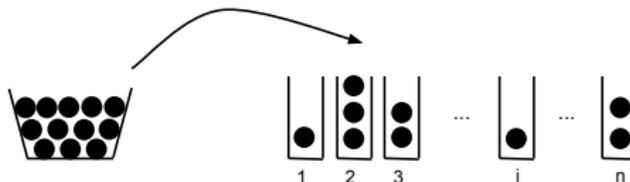
Begin

The bins are created based on their width (**data granularity**).



Round 1

The balls (energy consumption index values) from this specific period are distributed in the bins.



Probabilistic model - Game of balls and bins

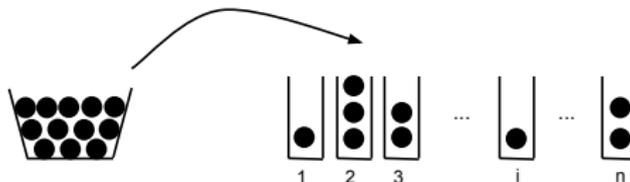
Begin

The bins are created based on their width (**data granularity**).

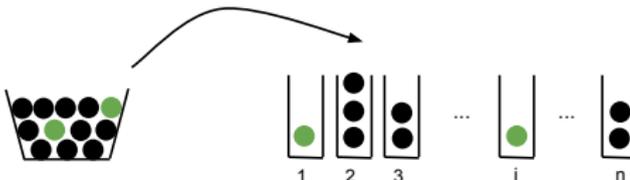


Round 1

The balls (energy consumption index values) from this specific period are distributed in the bins.



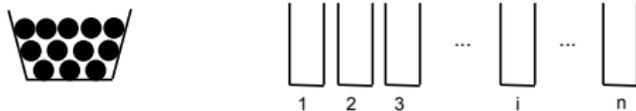
The balls that fall alone in their bin signify smart meters that are uniquely identified.



Probabilistic model - Game of balls and bins

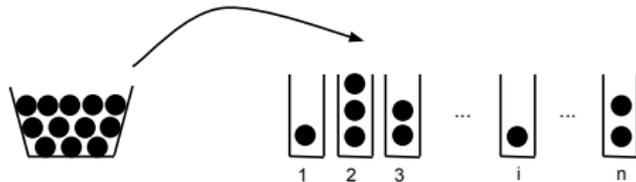
Begin

The bins are created based on their width (**data granularity**).

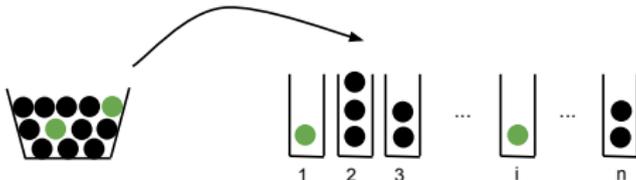


Round 1

The balls (energy consumption index values) from this specific period are distributed in the bins.

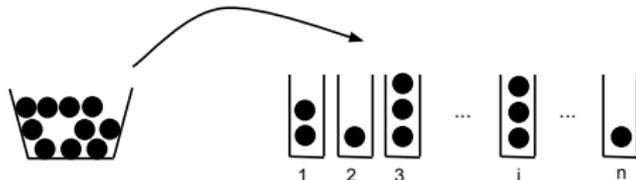


The balls that fall alone in their bin signify smart meters that are uniquely identified.



Round 2

The values for the smart meters identified last round are removed from this round data. The identification is repeated.



Assume a Poisson distribution of balls into bins²

- the expected number of consumption indexes identified uniquely at round j :

$$E_j[\text{bins with 1 ball}] = e^{-\frac{m_j \times w}{M}} \times m_j$$

For two consecutive rounds $j-1, j$

- at each round we remove the identified consumption indexes: $m_j = m_{j-1} - E_{j-1}[\text{bins with 1 ball}]$
- the number of consumption indexes at current round depends only on the number of consumption indexes at previous round and the number of bins considered: $m_j = m_{j-1} \times (1 - \exp(-\frac{m_{j-1} \times w}{M}))$

The game ends when

- all balls have been removed from the game: $m_j = 0$
- all the time periods with available data have been used: $j > T$

²Adapted from: M. Mitzenmacher and E. Upfal - *Probability and computing: Randomized algorithms and probabilistic analysis*, Cambridge University Press, 2005

Adversarial strategy modeled as a game of balls and bins

- We use the sizes of the different LF datasets as input for the **game of balls and bins**

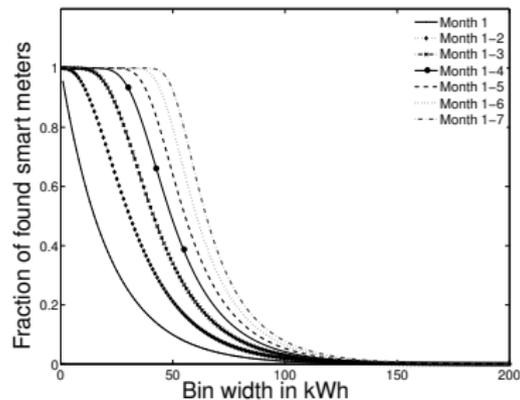
Actual execution of adversarial strategy

- We run the **adversarial strategy algorithm** on the different LF datasets

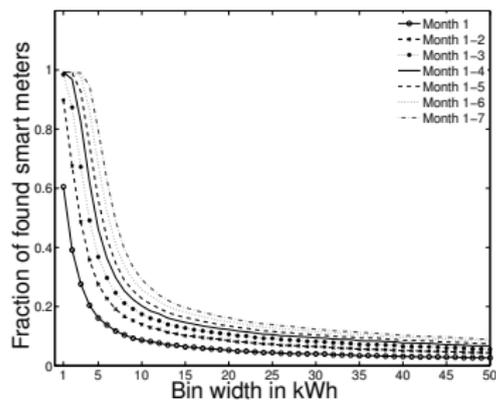
What is the effect of data **granularity** and data **timespan** on the ability of the adversary to identify a large number of customers?

Results - 7 months period

Bins and balls model



Real world data (RW data)



What is the effect of data **granularity** and data **timespan** on the ability of the adversary to identify a large number of customers?

Data granularity and data timespan

Granularity 1 kWh

Time period	Newly found smart meters		Total found smart meters %	
	Model	RW data	Model	RW data
m_1	18,461	11,698	95.4%	60.5%
m_2	871	5,655	99.9%	89.7%
m_3	2	1,669	100 %	98.3%
m_4	0	155	100 %	99.1%
m_5	0	11	100 %	99.2%
m_6	0	11	100 %	99.3%
m_7	0	10	100 %	99.3%
Total	19,334	19,209	100 %	99.3%

Granularity 10 kWh

Time period	Newly found smart meters		Total found smart meters %	
	Model	RW data	Model	RW data
m_1	12,182	1,670	63.0%	8.6%
m_2	6,029	1,027	94.1%	13.9%
m_3	1,093	671	99.8%	17.4%
m_4	30	543	100 %	20.2%
m_5	0	487	100 %	22.7%
m_6	0	579	100 %	25.7%
m_7	0	651	100 %	29.1%
Total	19,334	5,628	100 %	29.1%

- A change in the granularity of data reported monthly can significantly reduce the number of identified smart meters
- If laws and regulations allow → customers can opt for this type of reporting to gain extra privacy

Summary - de-anonymization

- Data granularity and data timespan have an important influence in AMI data privacy
- These two characteristics should be taken in consideration when releasing datasets to 3rd parties
- Even with the simple model a large number of smart meters can be identified uniquely based on their energy consumption

- The Advanced Metering Infrastructure (AMI)
- AMI data - utility and privacy
- Privacy issues - de-anonymization and de-pseudonymization
- **Differential-privacy and AMI data**
- AMI data application - energy load forecast using DP-aggregated AMI data
- Conclusion

Differential privacy³ (DP)

General Definition

A randomized function M gives ϵ -differential privacy for all data sets D and D' differing in at most 1 element, and all $S \subseteq \text{Range}(M)$, if

$$\Pr[M(D) \in S] \leq \exp(\epsilon) \times \Pr[M(D') \in S]$$

Noise addition

For $f : D \rightarrow R^d$, the mechanism M , which adds independently generated noise following the Laplace distribution $\mathcal{L}(\Delta f / \epsilon)$ to each of the d output terms, enjoys ϵ -differential privacy.

Mechanism's Sensitivity

For $f : D \rightarrow R^d$, the L_1 sensitivity of f is $\Delta f = \max_{D, D'} \|f(D) - f(D')\|_1$ for all D, D' differing in at most 1 element.

³Dwork, C., Naor, M., Pitassi, T. and Rothblum, G.N., 2010. Differential privacy under continual observation. In Proceedings of the forty-second ACM symposium on Theory of computing (pp. 715-724). ACM.

Differential privacy - maximizing the utility

For $f : D \rightarrow R^d$, the L_1 sensitivity of f is $\Delta f = \max_{D, D'} \|f(D) - f(D')\|_1$ for all D, D' differing in at most 1 element.

- for binary data ($\{0, 1\}$) the sensitivity is at most 1
- for real data (\mathcal{R}) the sensitivity is $\infty \rightarrow$ infinite noise, no utility

Differential privacy - maximizing the utility

For $f : D \rightarrow R^d$, the L_1 sensitivity of f is $\Delta f = \max_{D, D'} \|f(D) - f(D')\|_1$ for all D, D' differing in at most 1 element.

- for binary data ($\{0, 1\}$) the sensitivity is at most 1
- for real data (\mathcal{R}) the sensitivity is $\infty \rightarrow$ infinite noise, no utility

Solution?

Differential privacy - maximizing the utility

For $f : D \rightarrow R^d$, the L_1 sensitivity of f is $\Delta f = \max_{D, D'} \| f(D) - f(D') \|_1$ for all D, D' differing in at most 1 element.

- for binary data ($\{0, 1\}$) the sensitivity is at most 1
- for real data (\mathcal{R}) the sensitivity is $\infty \rightarrow$ infinite noise, no utility

Solution?

Limit the noise by bounding ^a
the sensitivity to a value B .

$$\mathcal{L}(\Delta f/\epsilon) \rightarrow \mathcal{L}(B/\epsilon)$$

^aGulisano, V., Tudor, V., Almgren, M. and Papatriantafilou, M., 2016, May. BES: Differentially Private and Distributed Event Aggregation in Advanced Metering Infrastructures. In Proceedings of the Second ACM International Workshop on Cyber-Physical System Security (pp. 59-69). ACM.

Simple DP aggregation - employing Δf

Simple Aggregation

$$S = \sum_{i=1}^n y_t^i, \text{ where } y_t^i \in [0, E].$$

Noise addition

$$\mathcal{L}(k\Delta f/\epsilon), \text{ where } k = \lceil WS/WA \rceil.$$

Measuring the error

$$\left| \frac{S - (S + \mathcal{L}(k\Delta f/\epsilon))}{S} \right| = \left| 0 - \underbrace{\frac{\mathcal{L}(k\Delta f/\epsilon)}{S}}_{Err_{noise}} \right|$$

Simple DP aggregation - employing Δf

Simple Aggregation

$$S = \sum_{i=1}^n y_t^i, \text{ where } y_t^i \in [0, E].$$

Noise addition

$$\mathcal{L}(k\Delta f/\epsilon), \text{ where } k = \lceil WS/WA \rceil.$$

Measuring the error

$$\left| \frac{S - (S + \mathcal{L}(k\Delta f/\epsilon))}{S} \right| = \left| 0 - \underbrace{\frac{\mathcal{L}(k\Delta f/\epsilon)}{S}}_{Err_{noise}} \right|$$

if $\Delta f \rightarrow \infty$ (very large consumption values) then the noise introduced by $\mathcal{L}(k\Delta f/\epsilon)$ will be very large

Bounded DP Aggregation - employing B

Bounded Aggregation

$$S_B = \sum_{i=1}^n \min(y_t^i, B), \text{ where } B \in [0, E].$$

Noise addition

$$\mathcal{L}(kB/\epsilon), \text{ where } k = \lceil WS/WA \rceil.$$

Measuring the error

$$\left| \frac{S - (S_B + \mathcal{L}(kB/\epsilon))}{S} \right| = \underbrace{\frac{S - S_B}{S}}_{Err_{approx}} - \underbrace{\frac{\mathcal{L}(kB/\epsilon)}{S}}_{Err_{noise}}$$

Bounded DP Aggregation - employing B

Bounded Aggregation

$$S_B = \sum_{i=1}^n \min(y_t^i, B), \text{ where } B \in [0, E].$$

Noise addition

$$\mathcal{L}(kB/\epsilon), \text{ where } k = \lceil WS/WA \rceil.$$

Measuring the error

$$\left| \frac{S - (S_B + \mathcal{L}(kB/\epsilon))}{S} \right| = \underbrace{\frac{S - S_B}{S}}_{Err_{approx}} - \underbrace{\frac{\mathcal{L}(kB/\epsilon)}{S}}_{Err_{noise}}$$

Limit the noise by bounding the sensitivity to B .
How to choose B ?

Bes - choosing a bound B

Using open data repositories

Compute B on data from an already public dataset or which can easily be made public.

Bes - choosing a bound B

Using open data repositories

Compute B on data from an already public dataset or which can easily be made public.

Use a differentially private mechanism to compute B

Choose the bound B among a set of candidate values $\mathcal{B} = \{B_1, \dots, B_o\}$ with a given mechanism M run over a dataset (D_{explore}) containing the events used to quantify the utility of each individual bound in \mathcal{B} .

Bes - choosing a bound B

Using open data repositories

Compute B on data from an already public dataset or which can easily be made public.

Use a differentially private mechanism to compute B

Choose the bound B among a set of candidate values $\mathcal{B} = \{B_1, \dots, B_o\}$ with a given mechanism M run over a dataset (D_{explore}) containing the events used to quantify the utility of each individual bound in \mathcal{B} .

Most Common B (MCB)

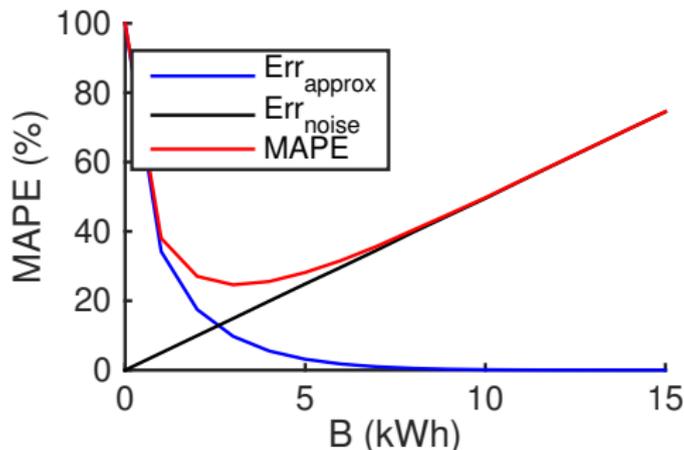
MCB aims at finding the bound $B \in \mathcal{B}$ resulting in the **minimum error for the majority of the SMs**.

High Enough B (HEB)

HEB looks for the bound $B \in \mathcal{B}$ for which **at least p of the n smart meters observe an error lower than the one observed for any higher bound B_j** .

Error composition - DP bounded aggregation

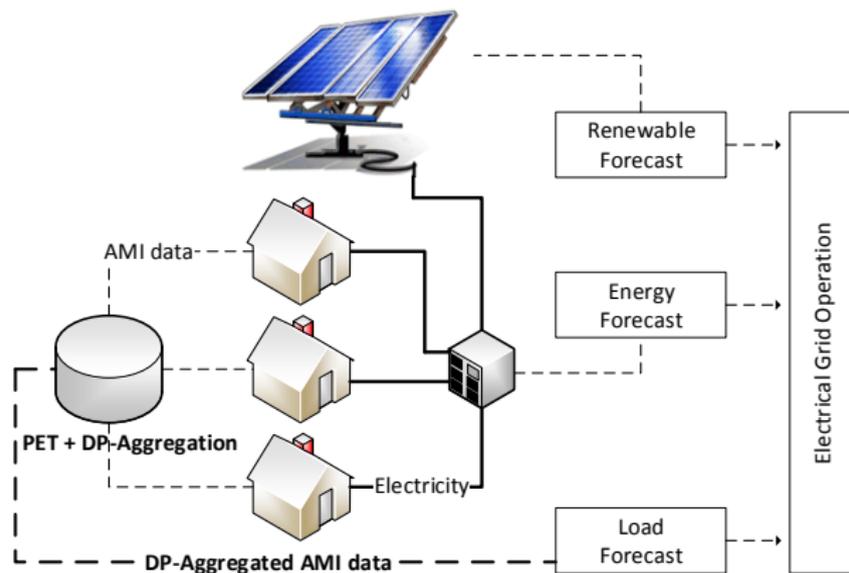
$$\left| \frac{S - (S_B + \mathcal{L}(kB/\epsilon))}{S} \right| = \underbrace{\frac{S - S_B}{S}}_{Err_{approx}} - \underbrace{\frac{\mathcal{L}(kB/\epsilon)}{S}}_{Err_{noise}}$$



for $B = \{0, 1, \dots, 15\}$ kWh.

- The Advanced Metering Infrastructure (AMI)
- AMI data - utility and privacy
- Privacy issues - de-anonymization and de-pseudonymization
- Differential-privacy and AMI data
- **AMI data application - energy load forecast using DP-aggregated AMI data**
- Conclusion

AMI load forecast scenario

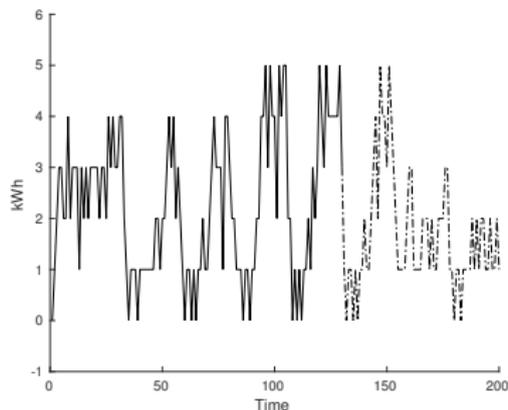


Using data for load forecast while protecting customers' privacy.

Load forecast

Load forecast

Predicting energy load based with the help of a model built on historical records.



Short term load forecast

Aims to predict consumption for short time frames, typically one hour to one week.

Load forecast - Model parameters

$$\hat{y} = \beta Y + \gamma W + \delta D + \alpha A$$

Y - energy load data,
 W - weather-related data,
 D - calendar-related data,
 A - anthropological data.

Variable	Characteristic	Forecast	Privacy
Y	Granularity	+	-
Y	Data collection periodicity (sampling)	-	+
Y	Dataset size (aggregated # customers)	+	+
Y	Training window size (duration)	+	-
Y	Test window size (duration)	+	-
Y	Predicted horizon (duration)	-	-
W	Temperature	+	neutral
D	Day of week	+	neutral*
A	Anthropological data	++	--

Persistence model (Seasonal Naïve) - PM

$$\hat{y}_t = y_{t-24}$$

Linear regression model 1 - LR1

$$\hat{y}_t = \beta_1 y_{t-24} + \beta_2 y_{t-48} + \beta_3 y_{t-72}$$

Linear regression model 2 - LR2

$$\hat{y}_t = \beta_1 y_{t-24} + \beta_2 y_{t-48} + \beta_3 y_{t-72} + \beta_4 \hat{T}_t + \beta_5 D(t)$$

LR1 and LR2 - Adapted from Y. Iwafune et al., "Short-term forecasting of residential building load for distributed energy management," in *Energy Conference (ENERGYCON), 2014 IEEE International*, May 2014, pp. 1197–1204.

Simple aggregation

$$y_t = \sum_{i=1}^N y_t^i$$

DP-aggregation

$$y_{t_{DP}} = \sum_{i=1}^N y_t^i + \mathcal{L}(B/\epsilon)$$

DP-aggregation deters de-anonymization attacks^a.

^aTudor, V., Almgren, M., and Papatriantafidou, M., 2015. A study on data de-pseudonymization in the smart grid. In Proceedings of the Eighth European Workshop on System Security (p. 2).

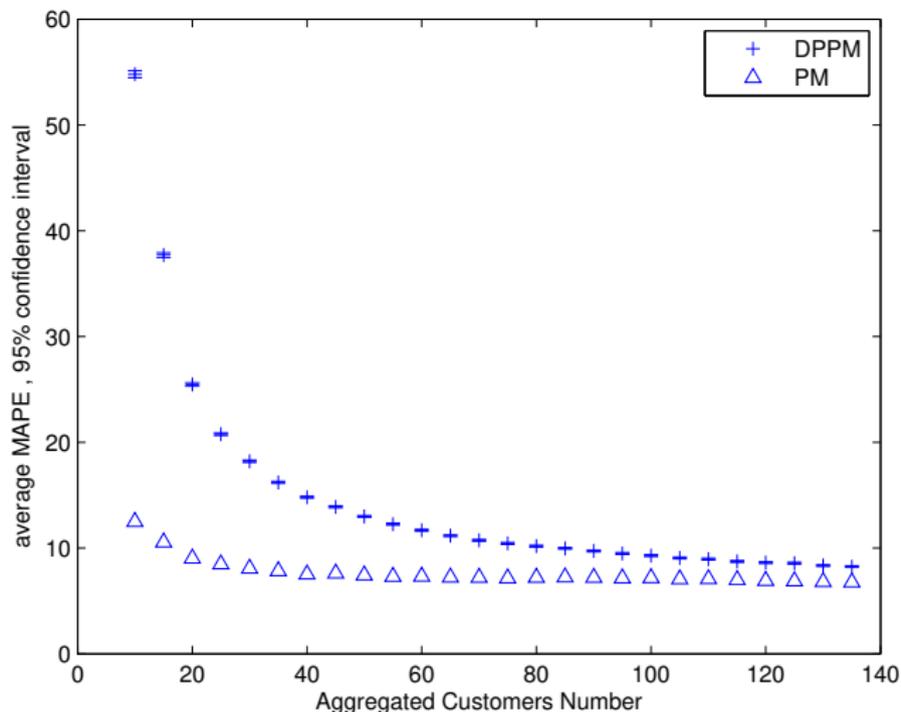
Evaluate the accuracy of three forecast models

- using simply aggregated AMI data
- using DP-aggregated AMI data
- on a variable number of customers

#	Name	Train Window	Forecast Horizon	Forecast Window	Number of Customers
1	PM	-	24h	1440h	10—135
2	DPPM	-	24h	1440h	10—135
3	LR1	1440h	24h	1440h	10—135
4	DPLR1	1440h	24h	1440h	10—135
5	LR2	1440h	24h	1440h	10—135
6	DPLR2	1440h	24h	1440h	10—135

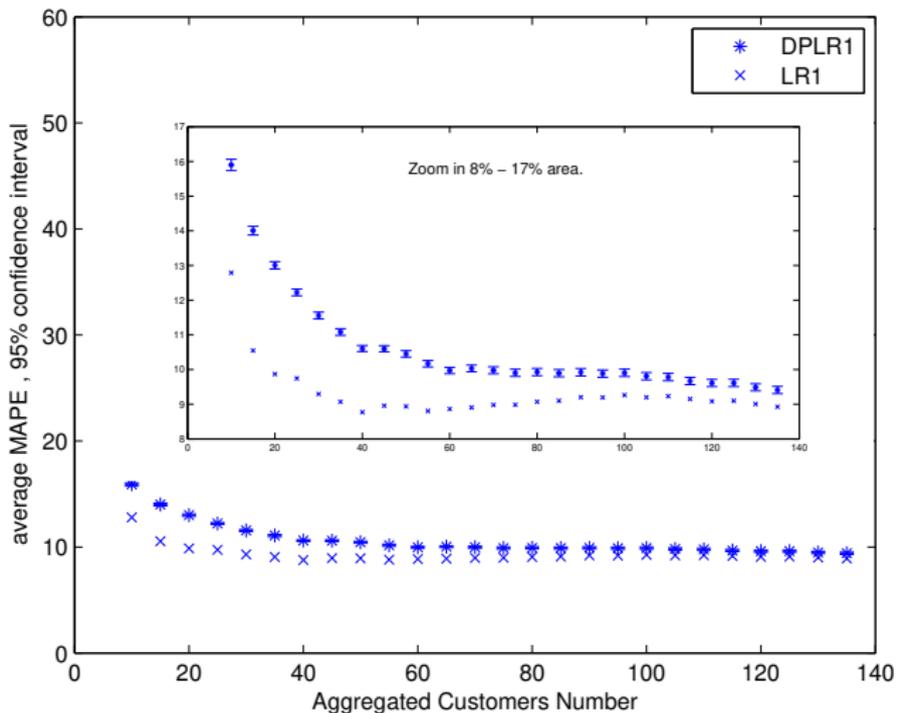
Average Mean Absolute Percentage Error (MAPE) for Persistent Method (PM)

24h forecast horizon, 100 tests/day and 60 predicted days



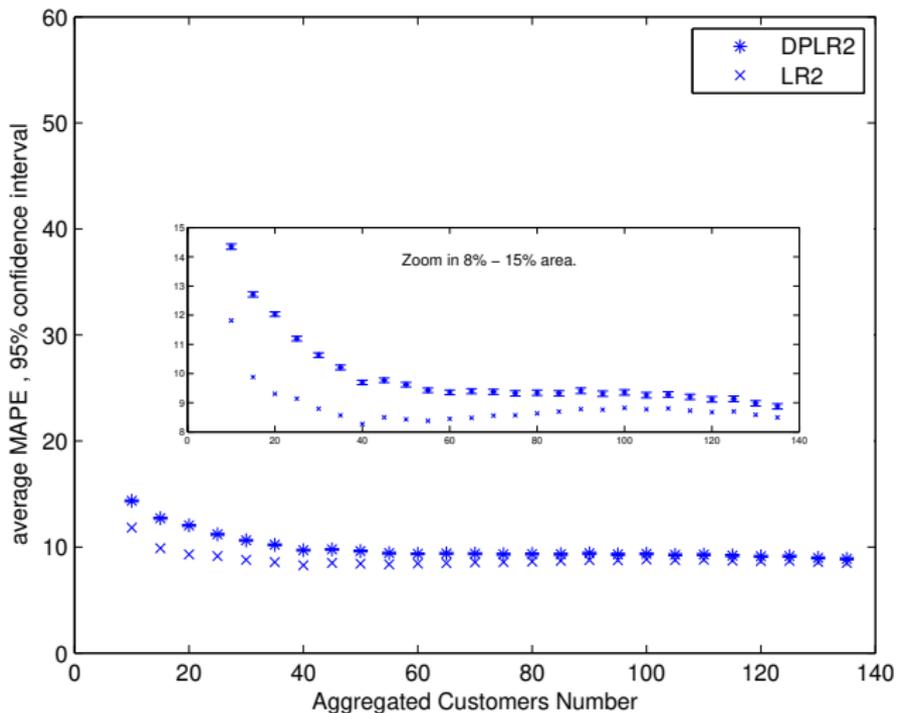
Average MAPE for Linear Regression model 1 (LR1)

24h forecast horizon, 100 tests/day and 60 predicted days



Average MAPE for Linear Regression model 2 (LR2)

24h forecast horizon, 100 tests/day and 60 predicted days



Which forecast method to use?

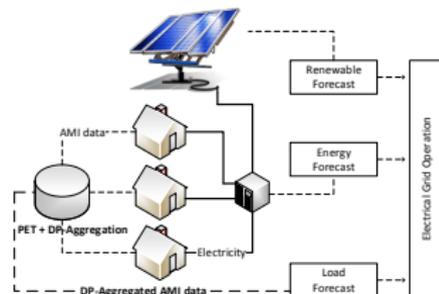
Number of Customers	PM	DPPM	LR1	DPLR1	LR2	DPLR2
10	12.50	54.80	12.79	15.90	11.82	14.35
35	7.84	16.22	9.07	11.08	8.57	10.21
50	7.43	13.00	8.94	10.45	8.43	9.62
100	7.17	9.30	9.26	9.91	8.82	9.35
135	6.78	8.25	8.92	9.41	8.49	8.87

Average MAPE (60 predicted days, 100 tests/day for DP methods)

$$\text{PM: } \hat{y}_t = y_{t-24}$$

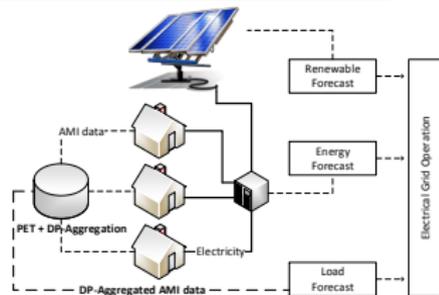
$$\text{LR1: } \hat{y}_t = \beta_1 y_{t-24} + \beta_2 y_{t-48} + \beta_3 y_{t-72}$$

$$\text{LR2: } \hat{y}_t = \beta_1 y_{t-24} + \beta_2 y_{t-48} + \beta_3 y_{t-72} + \beta_4 \hat{T}_t + \beta_5 D(t)$$



Summary - DP data application

- Differential privacy can successfully be employed in AMI data applications such as energy consumption forecasting.
- Short term forecast methods employing small-scale AMI data perform well and they can receive a boost in accuracy by further integrating **privacy neutral information**.
- Compared with their classical counterparts, the noise added by the forecast methods utilizing DP-aggregated data will introduce a **small prediction error**. This error increases with a decrease in the customers' group size.



That's all folks!

- The Advanced Metering Infrastructure (AMI)
- AMI data - utility and privacy
- Privacy issues - de-anonymization and de-pseudonymization
- Differential-privacy and AMI data
- AMI data application - energy load forecast using DP-aggregated AMI data

- Analysis of the impact of data granularity on privacy for the smart grid - V Tudor, M Almgren, M Papatriantafilou - Proceedings of the 12th ACM WPES 2013
- A study on data de-pseudonymization in the smart grid - V Tudor, M Almgren, M Papatriantafilou - Proceedings of Eurosec 2015
- BES: Differentially Private and Distributed Event Aggregation in Advanced Metering Infrastructures - V Gulisano, V Tudor, M Almgren, M Papatriantafilou - Proceedings of the 2nd ACM TCPS, 2016