# Blockchain Anonymity of Transactions and Household Location in a Smart Grid

## ANDREW MATHER

*Science and Enginering Faculty, Queensland University of Technology, 2 George St. Brisbane City, QLD, 4000, AUS*

*June 7, 2020*

---

This project aims to contribute to the study of user anonymity in blockchain for the Internet of Things. The research will explore methods to apply machine learning to deanonymise users on a smart grid with blockchain. Data stored on a decentralised blockchain is permanent and user privacy can be at significant risk. Research has found anonymity concerns in blockchain but IoT and smart grid contexts warrant further research.

The project analysis will include stored blockchain transactions and off-the-chain weather data. Section 1 will investigate the research topic and surrounding literature. Followed by covering the methodology to complete the project. The aim is to determine how effective machine learning is to deanonymise smart grid blockchain users. This should highlight long-term anonymity risks of blockchain by analysing permanent transaction data. The approach will first, source energy grid data and construct appropriate blockchain ledgers. Second, apply machine learning analysis on these blockchains to identify customer and their locations from past data. Third, incorporate off-the-chain weather data as customers also produce solar energy. Last, the project will investigate and test obfuscation methods to improve user privacy. Project outcomes will include a better understanding of user privacy risks in the smart grid scenario but also methods to mitigate these risks.

The progress reported so far is an initial completion of creating blockchain ledgers and applying transaction classification. A classification accuracy of 86.66% has been achieved with a multilayer perceptron approach when user's take no steps at protecting their privacy.

---

# CONTENTS

# 1. INTRODUCTION

## 1.1. Introduction

The massive growth in the Internet of Things (IoT) to collect, process, and send data via the Internet, requires a framework to handle all these devices. The IoT plays a role in many applications, for example, 'smart' devices in households, sensors, energy grids, and smart cities. Centralised IoT systems are challenged by cost, efficiency, and security as their size grows. A decentralised approach to the IoT's considerable mass of information [1] will become required, but user privacy and security challenges should be addressed.

Blockchain can handle this data as a decentralised ledger to record transactions carried out in a network. This is a developing field with wide potential uses, applications include cryptocurrency, financial systems, smart contracts, and non-monetary areas such as IoT and smart grids. Blockchain creates a level of anonymity for users through cryptographic means using private and public keys (PK).

The level of user anonymity from a permanent ledger is not studied in-depth in an IoT setting, despite the growing use of blockchain for the IoT. Studies on blockchain reveal malicious nodes can compromise user anonymity by classifying transactions using machine learning (ML) [2, 3] and using off-chain data [3]. The project aims to contribute to the study of user anonymity in an IoT blockchain smart grid. It will explore methods to deanonymise users using ML to analyse transaction [4] and off-chain weather data [5].

Note: The following subsections have been changed to increase specificity in the aim, objectives, and outcomes.

## 1.2. Thesis Statement

To investigate user anonymity in an IoT blockchain smart grid. The research will use machine learning to classify users and their location from blockchain transactions, and off-chain weather data. The purpose is to link transactions to users and identify their location to deanonymise them. Followed by providing techniques to enhance a user's privacy in the situation.

## 1.3. Context and Aim

Studies on blockchain reveal malicious nodes can compromise user anonymity. A node can link transactions and use off-chain data, e.g. weather data. Research has not studied user anonymity in IoT blockchain in-depth, despite widespread use. There is complexity in time series data and linked transactions, especially involving on and off-chain information.

The project aims to determine how effective ML is in deanonymising users in a smart grid implementing blockchain. This should highlight long-term anonymity risks of blockchain by analysing permanent transaction and historic weather data. The research is limited to a smart grid setting and ML classification as the analysis method. It is important also to determine and measure methods to improve user anonymity.

## 1.4. Objectives

Objective 1: Populate blockchain ledger from energy data.

- Source appropriate energy grid data.

- Convert sourced data into a blockchain format suitable for analysis.

Objective 2: Find the success rate of classifying blockchain transactions as specific users or locations.

- Establish the likelihood to extract a user's set of transactions.

- Investigate what classification methods are effective.

Objective 3: Find the success rate of linking user's transaction set to location with off-chain weather data.

- Measure the likelihood of locating a customer from their transactions and solar data.

- Investigate methods to compare customer and weather datasets.

Objective 4: Investigate the effectiveness of techniques to improve user anonymity.

- Investigate methods to increase user privacy in a smart grid blockchain with respect to objective 1 and 2.

- Measure the effect of several obfuscation techniques on transactions.

### 1.5. Significance and Outcomes

Blockchain for IoT has attracted tremendous attention recently. A huge volume of personalised data will become permanently stored in blockchains. Thus, it is critical to study the anonymity of the users in IoT. Identifying users and linking them to transactions in a smart grid, not only reveals private information, but also information such as when a home is unoccupied.

This research is undertaken to achieve:

- A better understanding of risks in adopting blockchain for smart grids.

- Objective 2 will establish the likelihood an attacker can link and extract a user's blockchain data.

- Objective 3 will establish the likelihood an attacker can combine a user's transaction data with weather information to deduce location.

- Objective 4 will analyse a range of privacy improving methods to measure their effectiveness.

The research is undertaken from an attacker's perspective on smart grid blockchain technology. The project involves trying to deanonymise others in the blockchain which a standard user would not attempt. The outcomes aim to benefit future users and ensuring their privacy in emerging technology.

## 2. BACKGROUND AND LITERATURE REVIEW

Literature from key areas of the project will be reviewed to highlight key concepts, locate information to aid the project's completion, and identify a research gap in section 2.6. First background information regarding blockchain and it's role in IoT and smart grids is covered. Then similar prior works in user anonymity and privacy are discussed. Last, machine learning techniques to classification users and transactions are covered. The related works section 2.4 and machine learning techniques 2.5 are extensions to the initial review.

### 2.1. Blockchain

Blockchain is a framework to create a public and universal distributed ledger. Bitcoin introduced blockchain [6] as a transaction ledger to ensure auditability, immutability, and nonrepudiation. Blockchain implements a method to reach consensus between unreliable parties. Whereas a standard process has a trusted third party (TTP), like a bank, responsible for transaction security. Blockchain properties remove the need for TTPs. Blockchain aims to store ordered transactions (blocks) linked to a previous block. Blocks contain a header, with a unique ID, and information [1]. Each block header stores the preceding block's hash to establish the links.

Network participants managing a blockchain are nodes or miners. They collate transactions into blocks to append to the blockchain. Networks use consensus algorithms to maintain trust and agreement to add a block. For example, Bitcoin uses a Proof of Work algorithm [6], and Ethereum a Proof of Stake [7]. Transactions utilise encryption, hashes, and public key (PK) cryptography. Digital signatures encrypt a document hash, signed with a private key, and a PK proves who signed [1]. Blockchain participants create anonymity by using PKs to conceal their real identity. Changing PKs for each transaction, as in Bitcoin [6], can improve user anonymity.

## 2.2. Blockchain for the Internet of Things

The IoT is physical devices connected to the internet and use communication networks to process data [1]. It makes devices 'smart' and gain computation and communications capabilities. Many devices cause large data traffic [8] and future applications could reach billions of devices [1]. Challenges limit IoT devices, including computing power, storage, and bandwidth. The data generated can cause bottlenecks. Blockchain can change how IoT networks operate with a decentralised framework [9, 10]. IoT networks could enjoy blockchain's lower costs, decentralised management, and inherent privacy [8]. A combination of IoT and blockchain looks to solve the challenges faced in IoT networks [10]. There are many IoT applications in daily life, businesses, and society, shown in Figure 1. Intelligent power distribution, or smart grids are the interest of this project.
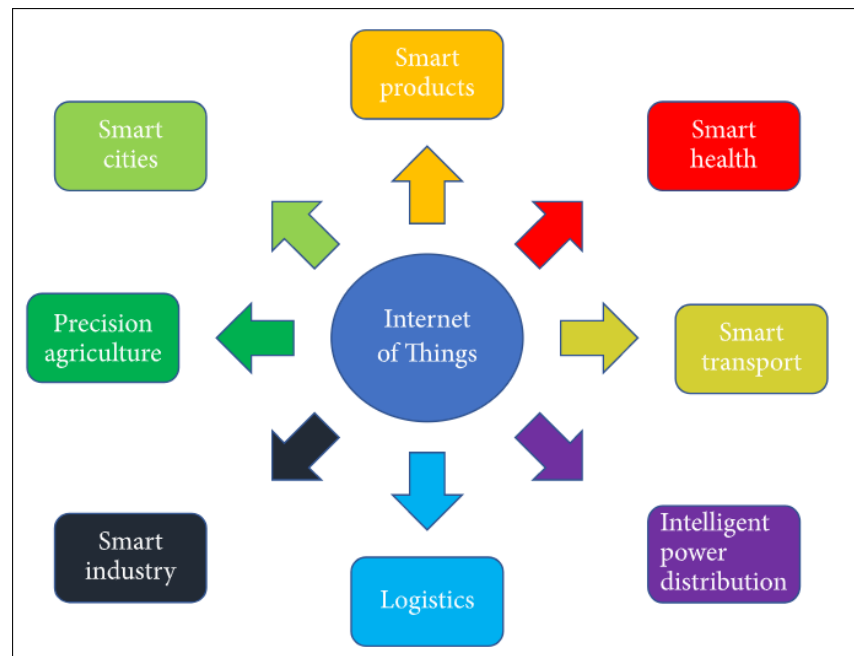


**Fig. 1.** IoT applications [1]

Energy systems and trading are developing quickly with benefits offered by the IoT [11]. Smart grids allow systems to communicate and optimise energy production, consumption [12], and thus utilisation [13]. They create the infrastructure to transfer energy between distributed producers and consumers. Some consumers can also be energy producers (prosumers) with solar energy [11]. There is continuously growing energy demand, supply, and quantities of energy sources. Making a desire for a decentralised energy system [14]. The authors at [11] highlight the challenges for such a system. First, managing transactions between users and the grid. Second, fluctuating supply from distributed renewable energy sources (e.g. rooftop solar). Third, TTPs lead to less efficient energy systems with more errors and operation costs. Blockchain is a likely path forwards to solve these issues in a smart grid [15].

Smart grids can decentralise behaviour using a blockchain framework. Blockchain and consensus algorithms can automate transactions and avoid TTPs. Other benefits can extend to real-time trading, anonymity features, and lower costs [16, 17]. The limits of IoT devices earlier are obstacles in implementing this system. The authors in [18] proposed a decentralised energy supply architecture. This provided on-demand energy for miners in an IoT network using microgrids.

## 2.3. Anonymity Concerns

Using IoT has many benefits but also increases exposure to new types of security, and privacy threats. IoT issues go beyond standard information and privacy concerns as it can relate to people's physical lives and security. Privacy relates to the large amounts of personal data used by smart devices [1]. Likewise, smart grids

will be complex networks, leading to privacy concerns [19] that need new approaches to solve [20]. Blockchain systems solve problems of centralised systems and increase resilience to failures and attacks [19]. This does not preclude new types of privacy risks. Blockchain anonymity research has been primarily into digital currencies [21, 22]. With its attractive properties, however, to create a trusted smart grid [19] further research is required in non-monetary IoT anonymity.

Privacy is the right a party has to disclose their information. Blockchain users create auditability through PKs while maintaining anonymity. The purpose is to mask a user's transactions, purchases, or information [1]. Relevant to a smart grid would be a participant's energy purchasing amounts, but also when, such as when a household is not consuming energy which implies the house is unoccupied. Privacy in blockchain should maintain transaction anonymity and have no ability to untie transactions [1]. Transaction anonymity means a transaction cannot be linked to a user, this is where different PKs are relevant. Untying transactions means transactions are not bound to user identities after routed through the network.

## 2.4.  Related Works - Extended

We can apply ML approaches to a blockchain to investigate if user transactions are linkable. As noted earlier, blockchain users change PKs each transaction to achieve anonymity. With supervised ML as suggested by [2], a malicious node could deanonymise a user by classifying transactions. One would use metrics such as the flow of inputs and outputs, to link users. An attacker can attempt deanonymisation with real-time network traffic, but this research will focus on blockchain and weather stored data. IoT networks are subject to privacy risks around the exposure of user activity patterns from sensed data [2].

The authors of [2] concluded cryptocurrency studies show users can be deanonymised from transaction patterns stored on a blockchain. Their research analysed blockchain transactions in an IoT and smart home environment with ML algorithms to classify devices. Analysis was performed as an informed and blind attacker on devices within the smart home. The attacker's aim was to identify what transactions belonged to what type of smart device. For example, identifying which transactions belong to a smart lock, thus inferring when an owner leaves their home. They populated a blockchain from real-world smart home network traffic. The attack method monitored the frequency of device transactions, using ML algorithms to compare with known frequency patterns of potential devices. Results showed an informed attacker being up to 90% accurate, and a blind attacker around 30%. This indicates a serious risk in privacy of devices using the blockchain. [2] also proposed methods to improve user privacy in the IoT blockchain. Three timestamp obfuscation methods reduced successful device classification by up to 30%. The techniques were; combining multiple packets into one transaction, merging ledgers, and adding random transactions delays. We can draw parallels between [2] and this project where smart devices become smart homes and we investigate the pattern of energy over time, as opposed to frequency.

Authors at [22] suggest an attacker in a multiple PK scenario needs to create a one-to-many mapping between users and addresses. The analysis process suggested by [22] involves three stages. First, the flow of blockchain transactions (nodes) in a transaction graph where PKs are the inputs and outputs. Second, the flow of payments between PKs in an address graph created from the transaction graph. Last, a user graph with the users and each PK that may belong to the same user; drawn from the previous information and blockchain heuristics. [23] used a full blockchain analysis to link users to public addresses.

For ML with unsupervised approaches, [3] is an example of clustering blockchain addresses. Clustering is a valuable method for ML problems [24], related to splitting data into groups. [3] takes a clustering approach to blockchain transactions and also off-the-chain data. Their scenario showed successful clustering of information, with off-the-chain data improving the accuracy. Clustering household energy patterns may show similarities between households located nearby.

[25] developed a method to deanonymise blockchain transactions using supervised machine learning to predict new entities. They perform multi-class classification to categorise a cluster of transactions. They use decision trees with random forest and gradient boosting algorithms. This paper has parallels with the project in attempting to categorise sets of transactions as belonging to different households in a smart grid.

The authors of [26] desired an accurate approach to compare weather data and power generation. They introduced linear and nonlinear time models for solar intensity prediction. This method could be a useful technique for the project's stage that compares solar data to a user's energy production.

A work relevant for the obfuscation techniques part of the project is [27]. The authors propose a privacy-preserving and data aggregation scheme. This will be useful for potential obfuscation techniques as they discuss dividing users into separate blockchains (ledgers) and using multiple pseudonyms (public keys) to protect a user's identity. There are similarities in the nature of these methods as [2].

## 2.5. Time Series Classification - Added

Time series data occurs in most tasks with human cognitive process. Classification problems with data that can be ordered, can be treated as a time series classification (TSC) problem [28]. Researchers have investigated many methods to effectively classify time series data [29]. Popular is nearest neighbour classifiers with a distance function if appropriate [30]. [30] also shows an ensemble of classifier's outperforms the individual components. These approaches use either an ensemble of decision trees (random forest) [31] or an ensemble of different types of discriminant classifiers [32].

[28] gives an overview of potential deep learning applications for TSC. The authors found for univariate and multivariate data, the top three types of networks were residual (ResNet), fully convolutional (FCN), and mulitlayer perceptron (MLP) networks. The MLP is traditional form of deep neural networks and was proposed in [33] as a baseline architecture for TSC.

Random forest is a decision tree machine learning approach used by [34] and [25] for TSC. The authors of [34] compared the effectiveness of different decision tree approaches for TSC. They found support vector machines performed poorly and ensembles are favoured when after optimal accuracy. The better performing ensembles were MultiBoost and AdaBoost.M1. However random forest performed similarly well and would be more favoured on larger datasets. This project will look to use a decision tree for classifying user blockchain transactions and will consider these options.

Cointegration is a potential statistical approach to compare time series. This is relevant for the project's comparison of household energy transactions to solar data. The authors in [35] investigated how to optimise wireless sensor networks in environmental monitoring. They used a statistical approach to cointegrate sets of time series data to select the optimal number of sensors. This was successful showing only 25% of the original sensors were not cointegrated. In particular [35] describes an analytical framework to analyse multivariate time series which relates to this project comparing solar data to user transactions.

## 2.6. Research Gap Identified

Research in blockchain user anonymity is developing and uses both transaction and off-the-chain analysis. Despite widespread usage of blockchain in IoT, user anonymity level is not yet studied thoroughly. Transactions involving data adds complexity, and it can be stored on or off-chain. The literature investigated shows research into IoT implementations of blockchain and some into the privacy, usually Bitcoin focused. The combination of machine learning analysis on stored data with the smart grid field is a new contribution. Privacy concerns abound as smart grids develop. It is important to understand the risks before storing user data on a permanent and public ledger.

## 3. METHODS AND PLAN

Additional specificity has been added to this section compared to the initial methodology. The timeline in section 3.6 has been updated to reflect the state of current progress.

### 3.1. Research Process

The following contains the research process broken down into steps. It covers the tasks required to address the thesis statement and objectives. Figure 2 describes an overview of the main phases planned for the project.

1. Background information and literature review on relevant papers for:

    (a) Blockchain-based IoT, smart grid and energy trading.

    (b) Privacy concerns in (a) and machine learning techniques to deanonymise users.

2. Source an appropriate energy dataset. It should contain a reasonable number of customers, energy use and solar production, and information that distinguishes users in different locations.

3. Convert energy data into blockchain ledgers. Include the ability to adjust the frequency of transactions, number of ledgers, and number of public keys per customer. Step 3 and 4 will address the first objective.

4. Perform machine learning analysis on blockchain transaction data to cover the second objective.

    (a) Test a variety of classification analysis approaches.

    (b) Apply to classifying transactions by customer or location.

    (c) Explore the effectiveness of unsupervised learning approaches briefly.

5. Perform machine learning on blockchain with off-the-chain weather data.

    (a) Investigate approaches to combine a user data set with solar data to reveal household location.

    (b) This is expected to include statistical analysis between separate datasets.

6. Suggest and evaluate methods to improve user anonymity.

    (a) Measure the change in privacy as consumers increase or decrease public keys used.

    (b) Measure the change in privacy as less ledgers than consumers are used.

    (c) Measure the effectiveness of timestamp obfuscation and combining transactions.

    (d) Combine the three methods to find combined benefit.
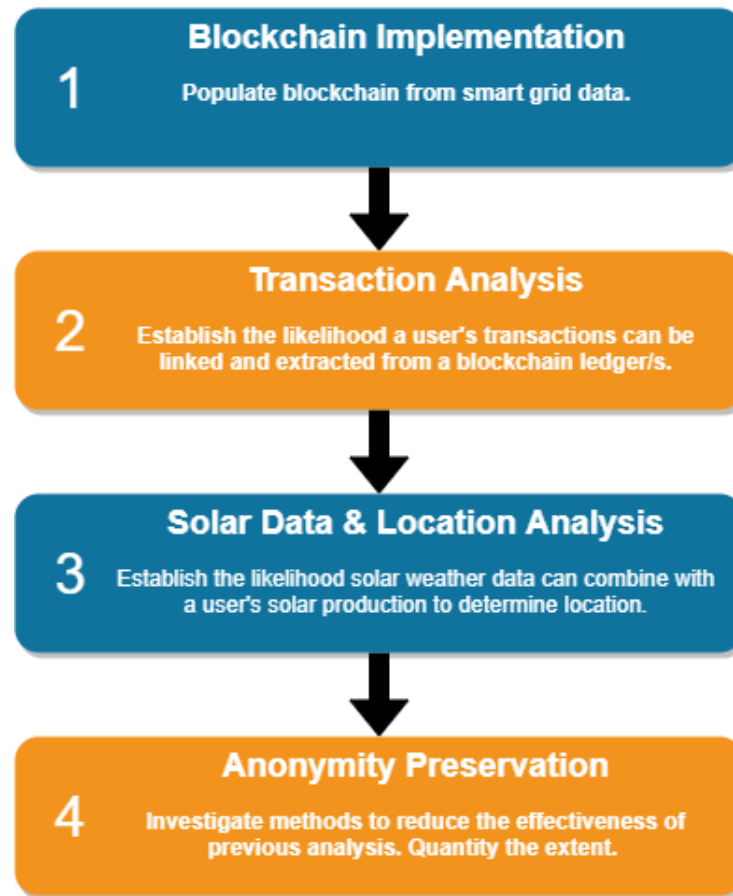
7. Deliver progress and final project reports.

**Fig. 2.** Phase design

### 3.2. Technical Frameworks

- Machine learning techniques:

  - Supervised learning: Classification, decision trees, neural networks.

  - Unsupervised learning: Clustering, correlation, cointegration.

- Python:

  - Populate blockchain from the data set/s.

  - Machine learning library: Scikit-learn.

  - Other key libraries: NumPy, MatplotLib, Pandas.

### 3.3. Data Collection

The project will describe the energy grid data and how to populate a blockchain implementation. Past energy use and generation data is sourced from Ausgrid solar home electricity data [4]. It contains Australian household data where households use and produce solar energy for the grid. The data set has half-hour data from 1 July 2010 until 30 June 2013 for 300 households across New South Wales. It includes energy consumption (on and off-peak) and generation. All customers are a full data set and Ausgrid performed quality checking. There is monthly data, but this makes less sense for a smart grid scenario.

Other datasets found such as [36] include less households, only 11 in this case, making anonymity hard to study. And are not situated in Australia with easy access off-the-chain data that Ausgrid and BOM provide.

The blockchain will be populated with transactions corresponding to the energy use of households in the dataset. Will assume each half-hour period is a communication between a smart meter and the grid. Thus, the process will generate a transaction for each period per household. Assumptions for this process will include:

- Real-time network traffic will be abstracted out. The focus is on attackers with access to permanently stored transaction information.

- Blockchain algorithms such as consensus algorithm are not required and will assume to pass each transaction. Patterns of transactions are not reliant upon these blockchain steps.

A single node will act as a miner collecting all transactions until a blocksize is reached. Then the miner creates a new block appended to the blockchain following the half-hour periods as the dataset. The results will be analysed with ML algorithms to deanonymise households to achieve objective two.

Relevant off-the-chain data is required for the second research objective. Ausgrid data provides postcodes for each household and therefore historical data from the Bureau of Meteorology will be easy to source and relevant. BOM data available at [5] can and solar exposure information for this project.

Other off-the-chain data worth investigating is provided by Ausgrid which ties in well with the original data. They provide average energy use by distribution zones at [37] and past outage data at [38]. The first provides general trends by region identifying areas of heavy use from little. The latter can explain anomalies in the data as a result of outages and prevent confusion in the learning models.

### 3.4. Data Analysis

Each phase of the project will perform a similar ML analysis on the acquired data. The project will use Python frameworks to perform supervised and unsupervised ML. Focusing on classification techniques. Python and the required libraries are freely available and appropriate to the situation.

The method will have an 'attacker' training machine learning models locally and measuring the ability of these models to predict users and their location. The next stage will use further classification and statistical approaches to link a user to the most similar solar data set. The last stage investigates obfuscation methods to increase resilience against ML attacks used. Likely suitable will be timestamp obfuscation methods, and combinations of different types of input and output data.

### 3.5. Project Management

To prepare to complete this project, the key research questions, objectives, and outcomes were detailed to highlight the focus of the project. A literature review was completed to reinforce the project direction and goals. This was extended to develop reasoning behind the selected classification methods implemented. A break-down of the phases of the project defines the necessary order of work in Figure 2 earlier. These processes will aid to ensure the project is completed in manageable sections and with appropriate quality control and testing. Each step shows the main research and investigation activities required.

### 3.6. Project Timeline

Figure 3 provides an overview of the project timeline. Objective two has been started earlier than initially planned by several weeks and the timeline has been amended to reflect this. This time was gained as the initial timeline overestimated the length of objective one. The change allows time to refine the initial transaction analysis after feedback and allocate greater time to the more important objectives three and four.
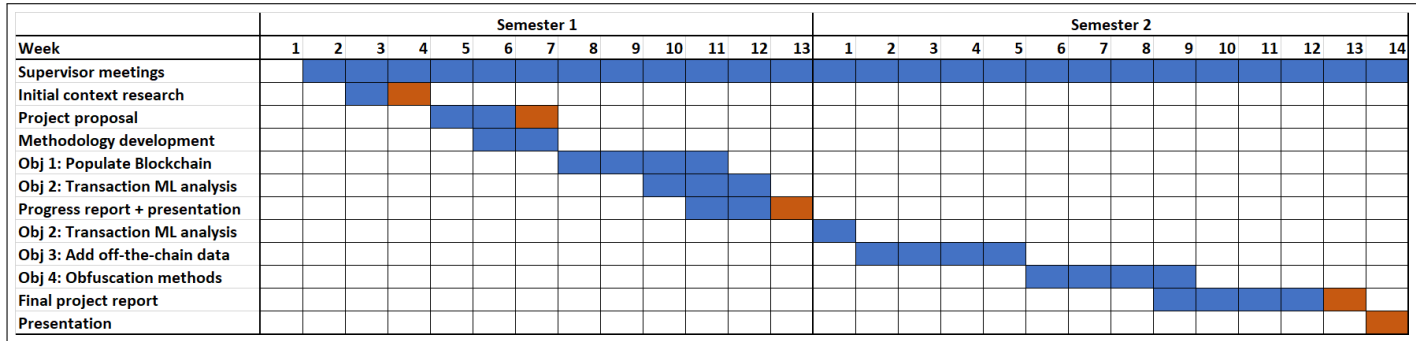
|  | Semester 1 | | | | | | | | | | | | | Semester 2 | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Week | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| Supervisor meetings | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| Initial context research | | | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | |
| Project proposal | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | |
| Methodology development | | | | | | ■ | ■ | | | | | | | | | | | | | | | | | | | | |
| Obj 1: Populate Blockchain | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | |
| Obj 2: Transaction ML analysis | | | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | |
| Progress report + presentation | | | | | | | | | | | | ■ | ■ | | | | | | | | | | | | | | |
| Obj 2: Transaction ML analysis | | | | | | | | | | | | | | ■ | | | | | | | | | | | | | |
| Obj 3: Add off-the-chain data | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | | |
| Obj 4: Obfuscation methods | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | | | | | |
| Final project report | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | |
| Presentation | | | | | | | | | | | | | | | | | | | | | | | | | | | ■ |

**Fig. 3.** Project timeline

# 4. RESEARCH PROGRESS

Progress has been made towards the project objectives (see section 1.4). Specifically, progress will be reported on objectives one and two. Objectives three and four will be completed in the project's second half. In additional to the progress discussed in this section, the literature review in sections 2.4 and 2.5 have been extended as research was required to make decisions on analytical techniques.

## 4.1. Energy Grid Data

The first step after project design and proposal was to source appropriate energy data. Ausgrid energy data [4] was found and evaluated to suit the project. The dataset has remained suitable for the project after progress so far. It is expected to contain all the information desired to complete the research. The following Table 1 shows a sample from the dataset.

**Table 1.** Example original energy data

| Customer | Generator | Postcode | Type | Date | 0:00 | 0:30 | 1:00 | ... | 23:30 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 3.78 | 2076 | CL | 01/07/2013 | 1.250 | 1.244 | 1.256 | ... | 1.081 |
| 1 | 3.78 | 2076 | GC | 01/07/2013 | 0.303 | 0.471 | 0.083 | ... | 0.068 |
| 1 | 3.78 | 2076 | GG | 01/07/2013 | 0 | 0 | 0 | ... | 0 |

This dataset contains the energy use and solar production of 300 households over three years. The data resolution is in half hour blocks which is nice from the perspective of creating blockchain transactions of a reasonable resolution. Each day a household's data is split into off-peak consumption (CL), general consumption (GC), and gross solar generation (GG). Additional features are the household generator size and postcode. Generator sizes range from 1kWh to 10kWH systems, but are primarily in the lower part of the range. This is a key driver of energy production and useful to have available. Key for the project is the postcode attribute allowing analysis to performed to classify household's from their data into postcodes using available solar data. Solar exposure data by area, in the time period of the energy data is available at [5].

## 4.2. Objective One - Populate Blockchain

Sourced energy data needs to first, be wrangled into a suitable format for time series classification, and second, have required features of a blockchain ledger added. This section describes the process to populate a blockchain from the energy data with various options. Visualisations of the dataset will be provided.

**Wrangle energy data**
The data was manipulated into a time series format with a datetime column removing the need for separate attribute columns for each time period. The purpose is so the blockchain ledgers have a row representing

a transaction with one timestamp and energy amount. This process drastically increases the size of the data which does create long analysis times for high time resolutions. An option added was to also produce datasets with different transaction time resolutions. These were, hourly, daily, and weekly transactions. Larger transaction sizes may more distinctly identify users, whereas high resolution data has lower values and more similarities. At this stage, the three years were left separate and 0 amount values kept. Table 2 shows a sample of the wrangled data.

**Table 2.** **Example wrangled energy data**

| Customer | Postcode | Type | Datetime | Amount |
|---|---|---|---|---|
| 1 | 2076 | CL | 01/07/2013 0:00 | 1.250 |
| 1 | 2076 | GC | 01/07/2013 0:00 | 0.303 |
| 1 | 2076 | GG | 01/07/2013 0:00 | 0 |
| 1 | 2076 | CL | 01/07/2013 0:30 | 1.244 |
| 1 | 2076 | GC | 01/07/2013 0:30 | 0.471 |
| 1 | 2076 | GG | 01/07/2013 0:30 | 0 |

**Create blockchain ledgers**

The energy data is now suitable for creating blockchain ledgers. There will be four transaction types:

- Genesis transaction $\rightarrow$ First transaction for a ledger or public key.

- On-peak consumption

- Off-peak consumption

- Solar production

Combining energy consumption and generation transactions into one per time period will be investigated as part of objective 3, mitigation and obfuscation strategies. Also investigated for this objective will be varying the number of public keys and blockchain ledgers. The ability to produce differing datasets to save time in the later project stages was implemented.

Adding several features is required to populate the blockchain. Each transaction has a hash (of its content) included as an identifier, and the previous transaction's hash to create a 'chain'. Also every household signs each transaction they generate with a PK. This produces the following structure of a transaction:
Hash | Previous Hash | Public Key | Timestamp | Transaction Type | Amount

For the second objective of transaction classification, two scenarios will be considered. First, an attacker's best case where each household has a separate ledger and one public key. Second, an attacker's worst case where all households are on a single ledger with new public keys for every transaction. Neither of these are realistic but will produce a bound of expectations. Realistic and privacy increasing (compared to best case) variations of ledgers and public keys will be analysed under objective four, obfuscation techniques. Households using a single public key is possible, but has a drastic reduction in privacy and is unlikely. On the other hand, new public keys require genesis transactions which cost the PK holder, cost-wise a new PK per transaction is unlikely.

Customer ID and postcode are left in the dataset for classifier training, but will be dropped from the test sets. Table 3 shows a sample of an output blockchain ledger.

**Table 3.** **Example blockchain ledger**

| Hash | PHash | PK | Customer | Postcode | Type | Datetime | Amount |
|---|---|---|---|---|---|---|---|
| Genesis | | $PK_1$ | 1 | 2076 | CL | 01/07/2013 0:00 | 1.250 |
| a | Genesis | $PK_1$ | 1 | 2076 | GC | 01/07/2013 0:00 | 0.303 |
| b | a | $PK_1$ | 1 | 2076 | GG | 01/07/2013 0:00 | 0 |
| c | b | $PK_1$ | 1 | 2076 | CL | 01/07/2013 0:30 | 1.244 |
| d | c | $PK_1$ | 1 | 2076 | GC | 01/07/2013 0:30 | 0.471 |
| e | d | $PK_1$ | 1 | 2076 | GG | 01/07/2013 0:30 | 0 |

### 4.3. Objective One - Blockchain Visualisation

It is important to understand trends in customer energy patterns. Figure 4 shows the pattern of energy use and generation of four customers over an example day. The figure legend includes the customer ID. Figure 4 brings two key insights. First, user consumption shows far more distinguishing features to separate users. Second, all customers have a similar solar generation trend (as expected from day/night cycles) but magnitudes have a large dependence on generator capacity. Customers shown, 37, 59, 102, 226 have solar capacities of 1.5, 2.8, 2.0, and 1.5kWh respectively. Displaying why customer 59's production is greater.

Figure 5 shows the consumption of the same consumers by day across the dataset. It is observed again consumption is a much better 'fingerprint' of a user's transactions. Trends, peaks, and troughs provide distinguishing features amongst these four consumers. The generation is also far more distinct in this view than Figure 4's single day.

Figure 6 shows the consumption of the same consumers by week across the dataset. This view seems quite similar to Figure 5 sharing many features, but with less fine detail variations and noise. While this looks visually to distinguish users easier the lost detail is important, as the results section 4.5 will show.
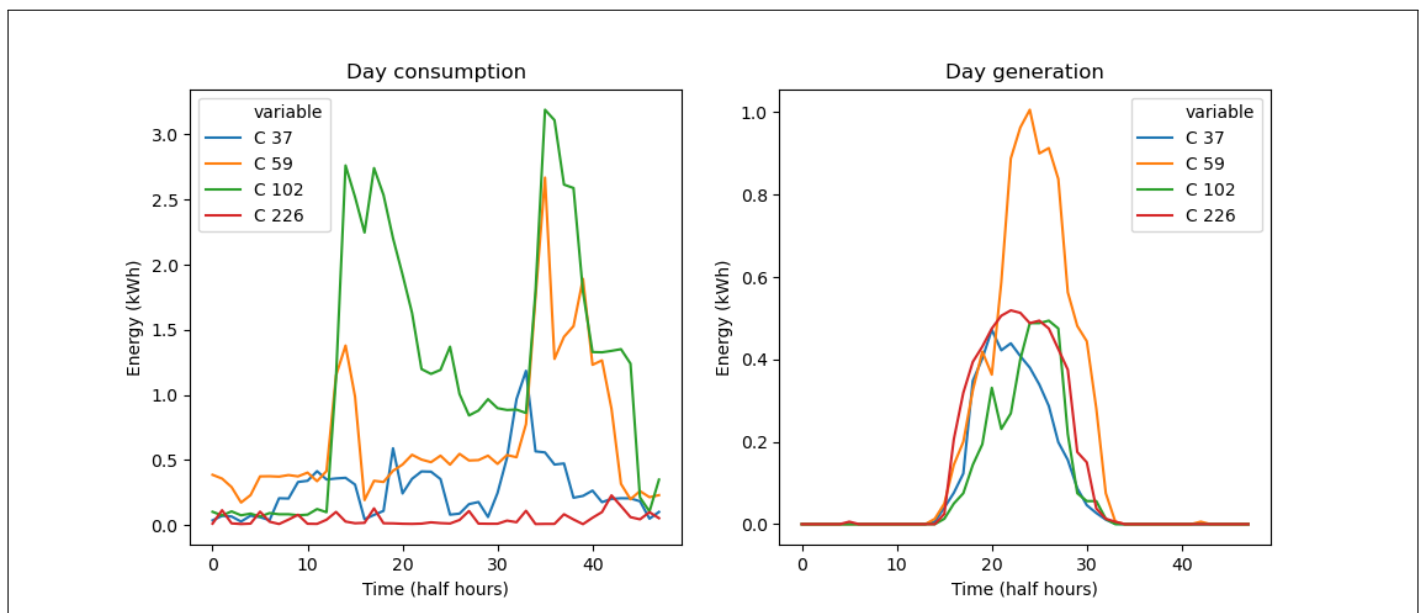


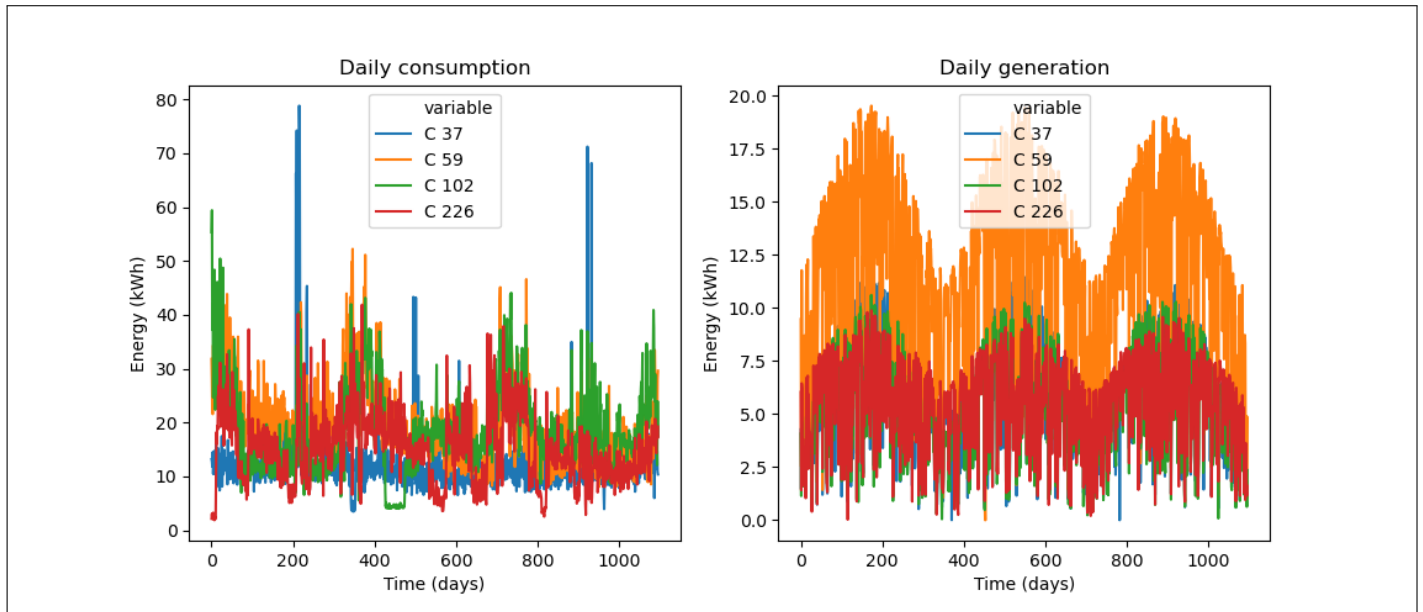**Fig. 4.** Energy pattern of four customers over a random single day

**Fig. 5.** Energy pattern of four customers by day over three years
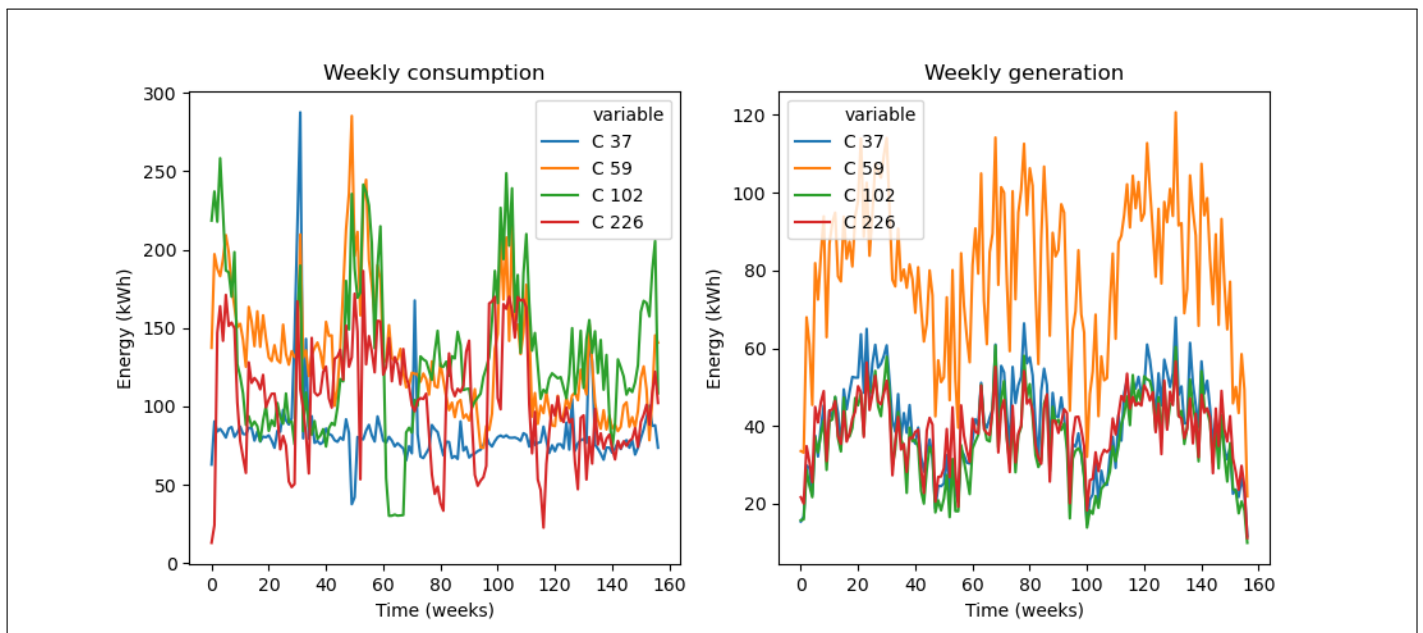


**Fig. 6.** Energy pattern of four customers by week over three years

## 4.4. Objective Two - Transaction Classification Methods

Attackers can construct a set of user's energy transactions by linking transactions emerging from the same PK or those with statistical similarity. This provides a time series which can fingerprint a user's energy consumption and generation. The generation is of more importance for objective three, and as shown, consumption is likely important to classify users in this stage. With set of user transactions, an attacker may be able to first, link ongoing transactions to the same user as an ongoing privacy risk, and second, potentially reveal household location with off-chain weather data.

The classification methods implemented are:

- K-Nearest Neighbours (KNN)

- Random Forest Decision Tree (RF)

- Multilayer Perceptron Neural Network (MLP)

**Selection of classifiers**
A KNN classifier was firstly selected as a simple baseline approach. The goal is to predict a category and the dataset has labelled data, thus classification is used over clustering. The dataset is not text based and thus a Stochastic Gradient Descent (SGD) or KNN classifier is suitable. Testing showed KNN achieving approximately double the accuracy of an SGD implementation on the same data.

Section 2.5 has discussed decision trees for time series classification (TSC). Random forest was an effective option, especially for large datasets like this project. Additionally, the random forest was suited to multivariable analysis and was easier to implement for the project with greater interest in objectives three and four.

Section 2.5 also covered suitable deep learning networks for TSC. While a residual network (a CNN) was reported best in the main paper discussed, MLP networks were also effective and may generalise larger multivariate data well. Part of this section involved implementing a recurrent neural network, LSTM. It however was found less effective for tabulated and time series data than both MLP and CNNs. Comparing the MLP results to a residual network will be listed as an improvement to this section.

**Training and test approach**

1. Preprocess data:

    (a) Categorical data made numerical and scaled.

    (b) Random train and test sets constructed. Customer, postcode, generator dropped from test sets.

    (c) Zero amount transactions stripped as these would not create blockchain transactions.

    (d) For worst case PKs are made unique. Ledgers combined so user transactions link in time order.

2. Run each classifier for weekly, daily, hourly, and half hourly transaction resolutions on:

    (a) Customer prediction: i) best case, ii) worst case

    (b) Postcode prediction: i) best case, ii) worst case

An initial attempt at clustering and cointegrating household solar production was made. This revealed clear clusters, unfortunately these align mainly by generator capacity and further research and analysis will be required.

## 4.5. Objective Two - Transaction Classification Results

This section presents the results of the classifiers and analysis approach discussed in section 4.4. Graphs of the results will be presented and discussed. Note, guesswork should expect average accuracies of 0.33% for customer and 1% for postcode. The best results were achieved by the MLP neural network classification, and highlights the importance of including a residual network, as suggested by literature, in ongoing work.

**K-Nearest Neighbours**
The first baseline classifier implemented was the KNN approach. Figure 7 displays the results of the KNN classification. The results showed a maximum accuracy of 14.08% for customer and 17.56% for postcode classification on hourly transaction resolution. This is an average and for example could be as large as 45% for predicting customer 138. Overall these results outperform guesswork but still quite poor. As expected because there are less options, accuracy in predicting postcode is greater than the customer.

Accuracy trends upwards as transaction time resolution increases, but drops for half hourly. This level of resolution perhaps masks consumer patterns with smaller differences harder to distinguish and noise more impactful. The worst case has a large impact on the KNN results, with customer classification dropping to at best, 2.22%.
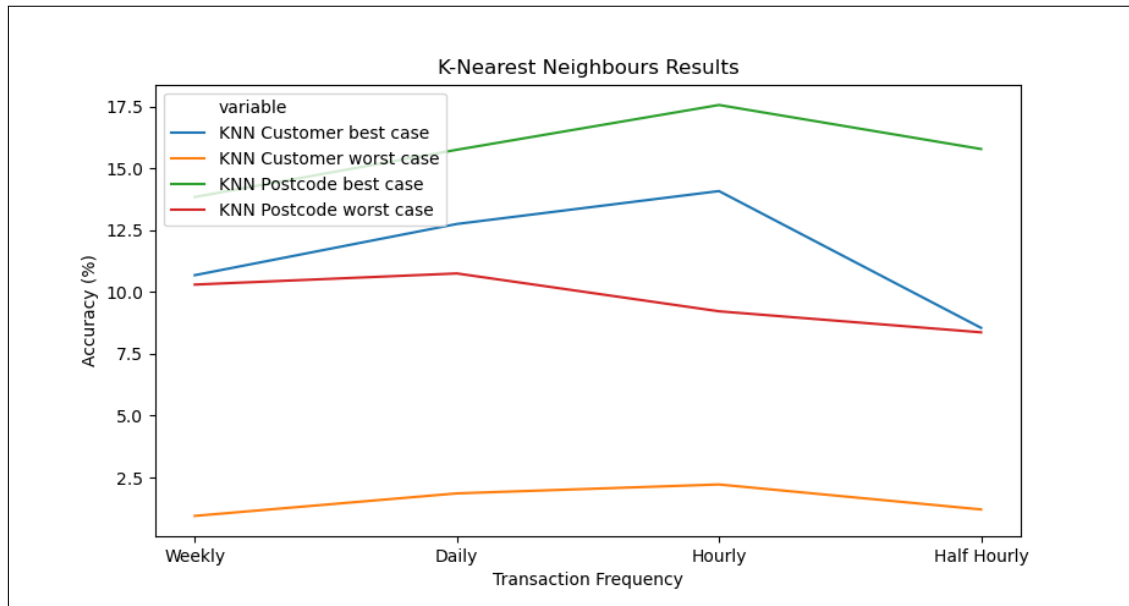


**Fig. 7.** Overall accuracy of KNN classification by transaction resolution and scenarios

**Random Forest**
Next a random forest decision tree is implemented and Figure 8 displays the classifier's results. The best case results are unusual at about 98% for both customer and postcode predictions across each time resolution. Meanwhile, the worst case was approximately guess work. The feature weights given by the model were almost entirely reliant on the public keys in best case. It was ensured the model was not being trained on public keys that matched the test set. These were changed so the model's couldn't simply 'learn' each user's public key. The outcome comes about when the decision tree recognises better than other classifiers a single user always uses the same public key. Therefore, after separating each user on a higher level, it slightly weights energy type and amount, it can immediately classify all its transactions correctly. While great, the classifier's inability on the worst case indicates realistic cases will experience a substantial and fast drop in performance.
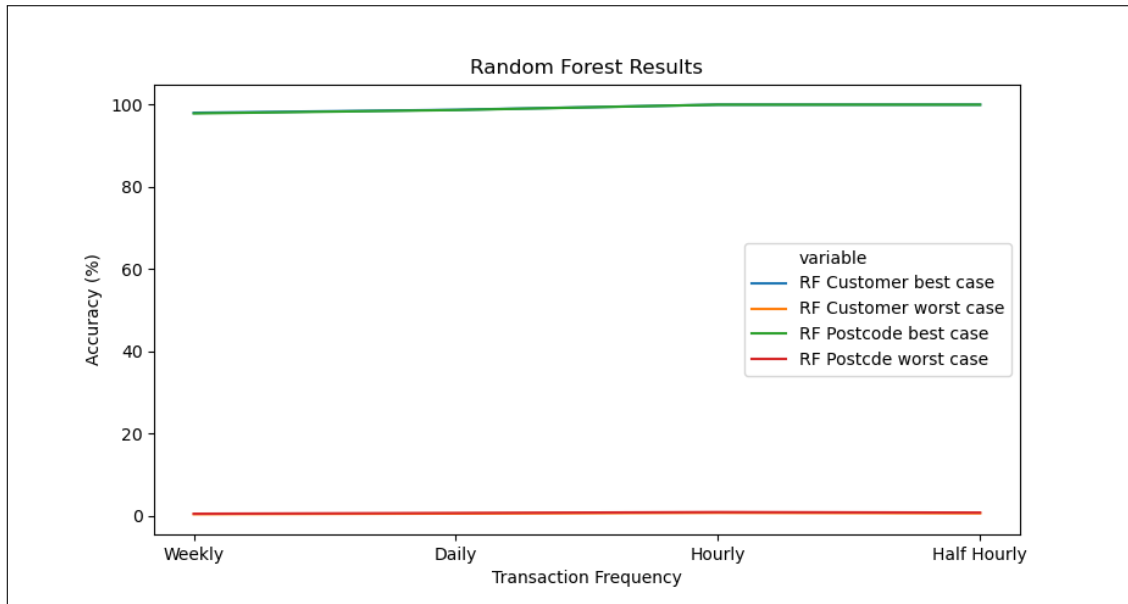
**Fig. 8.** Overall accuracy of RF classification by transaction resolution and scenarios

**Multilayer Perceptron**

The best classification approach used is the MLP network with results shown in Figure 9. The results showed a maximum accuracy of 86.66% for customer on half hourly resolution, and 45.46% for postcode classification on hourly resolution. This is an average and for example, 72 customers achieved 100% accuracy, and 10 0%. This large range of outcomes is shown by large standard deviations in the MLP results shown in Figure 10. Overall these results and the trend towards better accuracy with higher resolution are quite promising.

It is interesting accuracy for postcode prediction trends upwards, but drops for half hourly. However, this classifier appears to better handle the half hourly data than both KNN and RF approaches, perhaps less distracted by similarities when all transactions values are low and more similar.

MLP handles the worst case better than the previous models. The classifier can still achieve 2.28% and 11.48% accuracy for customer and postcode classification. This occurs on the weekly data where with a user's transactions no longer on separate chains, only the energy values can identify them. Thus the weekly resolution data has larger values which more distinctly separate households. This should indicate it will perform more effectively in the obfuscation cases that attempt to improve a user's privacy beyond the initial best case.

In general each classifier performs better on finer transaction resolution, with the simpler ones losing accuracy at the lowest half-hour granularity. MLP is the best performing model and is quite accurate. The best case accuracy of 86.66% and 72 households fully classified, shows an attacker chances at linking a user's transactions far above acceptable. Whether this can be used to deanonymise a user is subject to the project's next stage.
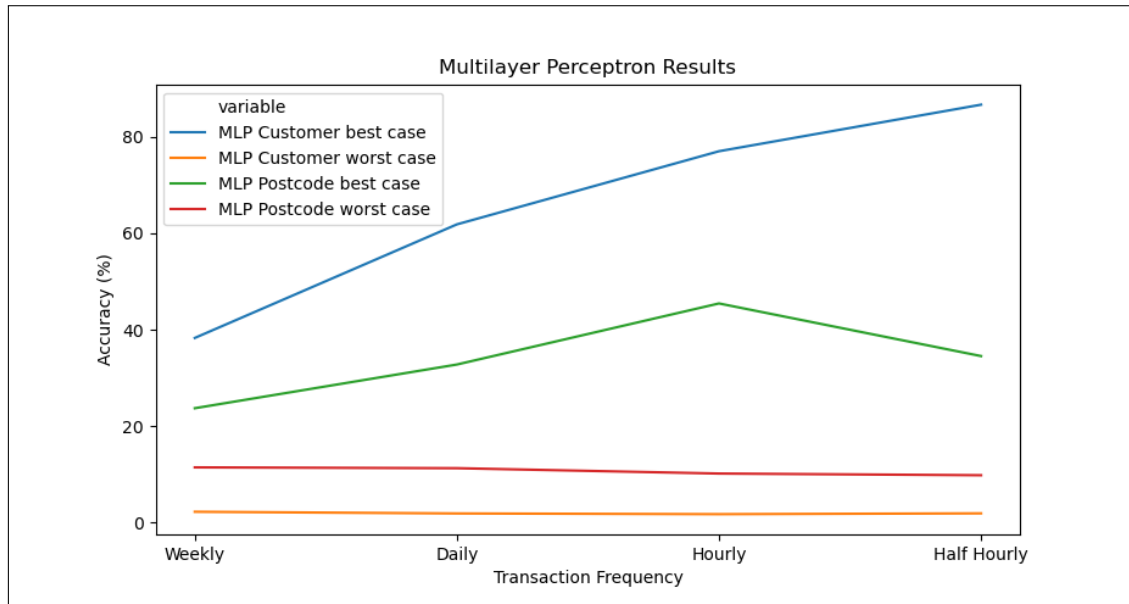
**Fig. 9.** Overall accuracy of MLP classification by transaction resolution and scenarios
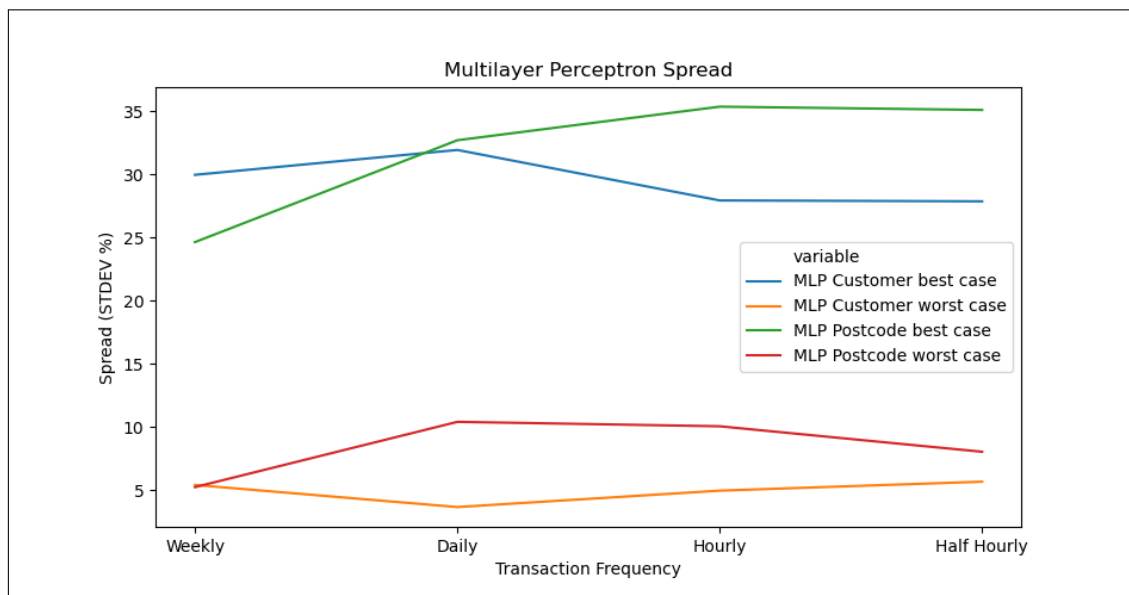


**Fig. 10.** Spread in standard deviation of the MLP classification results

**Improvements**

The following are features of objective two's progress suitable for improvement:

- Perform classification with a fully convolutional network (residual) and compare against the MLP used.

- Explore rank accuracy from the network classifier to indicate whether even if overall accuracy is achieved, an attacker can easily narrow the possible user's a set of PK transactions belong too.

- K-fold validation instead of randomised training and test sets.

## 5. NEXT STEPS AND CONCLUSION

The next stage of this project will implement feedback on the progress achieved to date, and perform a residual network classification analysis. Then remaining project stages are objectives three and four (see section 1.4). The research process from section 3.1 will be continued by measuring the likelihood a household's transaction energy production data can be linked to a postcode using solar exposure data. This will consist of several statistical or machine learning approaches which will require further research before implementing. An initial look at clustering energy production for objective three showed promise, if deeper links than generator capacity can be found. Objective three will look at the effectiveness of first clustering energy production. Initial attempt at this showed some promise. The fourth objective involves exploring several methods to improve user anonymity in the project's context. Such as, multiple public keys, ledgers, and timestamp obfuscation.

To achieve the progress so far I have had to learn about blockchain implementation and privacy concepts. But more importantly delve into machine learning approaches beyond standard classification and clustering concepts learnt previously. This field of data manipulation and analysis is massive and exploring it is interesting and a good learning experience.

The progress reported so far is an initial completion of creating blockchain ledgers and applying transaction classification. A classification accuracy of 86.66% has been achieved with a multilayer perceptron approach when if user's take no steps at protecting their privacy. This should highlight concerns for user privacy if no steps are taken to improve it.

# 6. REFERENCES

1. E. F. Jesus, V. R. L. Chicarino, C. V. N. de Albuquerque, and A. A. de A. Rocha, "A survey of how to use blockchain to secure internet of things and the stalker attack," Secur. Commun. Networks pp. 1–27 (2018).
2. A. Dorri, C. Roulin, R. Jurdak, and S. S. Kanhere, "On the activity privacy of blockchain for iot," (2019), pp. 258–261.
3. M. P. D. Ermilov and Y. Yanovich, "Automatic bitcoin address clustering," (2017), p. 461–466.
4. Ausgrid, "Solar home electricity data," https://www.ausgrid.com.au/Industry/Our-Research/Data-to-share/Solar-home-electricity-data (2014).
5. Bureau of Meteorology, "Climate data online," http://www.bom.gov.au/climate/data/ (2020).
6. S. Nakamoto, "Bitcoin: A peer-to-peer electronic cash system," (2008).
7. G. Wood, "Ethereum: A secure decentralised generalised transaction ledger," https://www.ausgrid.com.au/Industry/Our-Research/Data-to-share/Past-outage-data (2014).
8. M. Ferrag, M. Derdour, M. Mukherjee, A. Derhab, L. Maglaras, and H. Janicke, "Blockchain technologies for the internet of things: Research issues and challenges," IEEE Internet Things J. **6**, 2188–2204 (2019).
9. D. Puthal, N. Malik, S. Mohanty, E. Kougianos, and G. Das, "Everything you wanted to know about the blockchain: Its promise, components, processes, and problems," IEEE Consumer Electron. Mag. **7**, 6–14 (2018).
10. M. Swan, *Blockchain: Blueprint for a New Economy* (O'Reilly, Sebastopol, CA, USA, 2015), 1st ed.
11. T. Alladi, V. Chamola, J. Rodrigues, , and S. Kozlov, "Blockchain in smart grids: A review on different use cases," Sensors (Switzerland) **19** (2019).
12. R. Bayindir, I. Colak, G. Fulli, and K. Demirtas, "Smart grid technologies and applications," Renew. Sustain. Energy Rev. **66**, 499–516 (2016).
13. U. Ahsan and A. Bais, "Distributed big data management in smart grid," 26th Wirel. Opt. Commun. Conf. pp. 1–6 (2017).
14. L. Cheng, N. Qi, F. Zhang, H. Kong, and X. Huang, "Energy internet: Concept and practice exploration," (2017), p. 1–5.
15. M. Andoni, V. Robu, D. Flynn, S. Abram, D. Geach, D. Jenkins, P. McCallum, and A. Peacock, "Blockchain technology in the energy sector: A systematic review of challenges and opportunities," Renew. Sustain. Energy Rev. **100**, 143–174 (2019).
16. V. Hassija, G. Bansal, V. Chamola, V. Saxena, and B. Sikdar, "Blockcom: A blockchain based commerce model for smart communities using auction mechanism," (2019), p. 1–6.
17. G. Bansal, V. Hassija, V. Chamola, N. Kumar, and M. Guizani, "Smart stock exchange market: A secure predictive decentralised model," (2019), p. 1–6.
18. J. Li, Z. Zhou, J. Wu, J. Li, S. Mumtaz, X. Lin, H. Gacanin, and S. Alotaibi, "Decentralized on-demand energy supply for blockchain in internet of things: A microgrids approach," IEEE Transactions on Comput. Soc. Syst. **6**, 1395–1406 (2019).
19. B. Muhammad, J. Zhao, D. Niyato, L. Kwok-Yan, and X. Zhang, "Blockchain for future smart grid: A comprehensive survey," IEEE Transactions on Comput. Soc. Syst. (2019).
20. P. Kumar, Y. Lin, G. Bai, A. Paverd, J. S. Dong, and A. Martin, "Smart grid metering networks: A survey on security, privacy and open research issues," IEEE Commun. Surv. & Tutorials **21**, 2886–2927 (2019).
21. M. K. Khalilov and A. Levi, "A survey on anonymity and privacy in bitcoin-like digital cash systems," IEEE Commun. Surv. & Tutorials **20**, 2543–2585 (2018).
22. M. Conti, S. Kumar, C. Lal, and S. Ruj, "A survey on security and privacy issues of bitcoin," IEEE Commun. Surv. & Tutorials **20**, 2543–2585 (2018).
23. D. Ron and A. Shamir, "Quantitative analysis of the full bitcoin transaction graph," (2013), pp. 6–24.
24. Z. Ghahramani, "Unsupervised learning," Adv. lectures on machine learning **20**, 72–112 (2003).
25. H. H. S. Yin, K. Langenheldt, M. Harlev, R. R. Mukkamala, and R. Vatrapu, "Regulating cryptocurrencies: A supervised machine learning approach to de-anonymizing the bitcoin blockchain," J. Manag. Inf. Syst. **36**, 37–73 (2019).
26. Y. Wang, G. Cao, S. Mao, and R. M. Nelms, "Analysis of solar generation and weather data in smart grid with simultaneous inference of nonlinear time series," in *2015 IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS),* (2015), pp. 600–605.
27. Z. Guan, G. Si, X. Zhang, L. Wu, N. Guizani, X. Du, and Y. Ma, "Privacy-preserving and efficient aggregation based on blockchain for power grid communications in smart communities," IEEE Commun. Mag. **56**, 82–88 (2018).
28. H. I. Fawaz, G. Forestier, J. Weber, L. Idoumghar, and P.-A. Muller, "Deep learning for time series classification: a review," Data Min Knowl Disc **33**, 917–963 (2019).
29. B. A, L. J, B. A, L. J, and K. E, "The great time series classification bake off: a review and experimental evaluation of recent algorithmic advances," Data Min. Knowl. Discov. **31**, 606–660 (2017).
30. J. Lines and A. Bagnall, "Time series classification with ensembles of elastic distance measures," Data Min. Knowl. Discov. **29**, 565–592 (2015).
31. M. G. Baydogan, G. Runger, and E. Tuv, "A bag-of-features framework to classify time series," IEEE Transactions on Pattern Analysis Mach. Intell. **35**, 2796–2802 (2013).
32. A. Bagnall, J. Lines, J. Hills, and A. Bostrom, "Time-series classification with cote: The collective of transformation-based ensembles," IEEE Transactions on Knowl. Data Eng. **27**, 2522–2535 (2015).
33. Z. Wang, W. Yan, and T. Oates, "Time-series classification with cote: The collective of transformation-based ensembles," Int. joint conference on neural networks p. 1578–1585 (2017).
34. A. Jović, K. Brkić, and N. Bogunović, "Decision tree ensembles in biomedical time-series classification," **7476**, 917–963 (2012).
35. S. Bhandari, N. Bergmann, R. Jurdak, and B. Kusy, "Time series analysis for spatial node selection in environment monitoring sensor networks," Sensors **18**, 11–27 (2017).
36. Open Power System Data, "Household data," (2017).
37. Ausgrid, "Distribution zone substation data," https://www.ausgrid.com.au/Industry/Our-Research/Data-to-share/Distribution-zone-substation-data (2020).
38. Ausgrid, "Past outages," https://www.ausgrid.com.au/Industry/Our-Research/Data-to-share/Past-outage-data (2020).