

**FACULDADE DE ENGENHARIA DA UNIVERSIDADE DO PORTO**



**FEUP** FACULDADE DE ENGENHARIA  
UNIVERSIDADE DO PORTO

# **Cellular Time Activation Networks, a novel approach applied to photovoltaic anomaly detection**

**David da Silva Moreira Freire**

WORKING VERSION

Mestrado Integrado em Engenharia Eletrotécnica e de Computadores

Orientador: Cláudio Domingos Martins Monteiro

June 30, 2023

© David Freire, 2023

# **Resumo**

A proliferação de centrais fotovoltaicas de dimensão industrial levou à necessidade de métodos para detetar e classificar falhas nos seus componentes, sendo que estas que podem ter impactos económicos significativos. Neste trabalho, explora-se o estado da arte das ferramentas de deteção de falhas e estimação do estado aplicadas ao campo dos sistemas PV, com foco na compreensão do seu funcionamento, identificando-se pontos fortes e possíveis limitações. Relevar-se-á quais os métodos baseados em aprendizagem computacional mais utilizados. Ainda assim, reconhece-se o contributo dos diversos domínios para colmatar este tipo de problema, desde a teoria dos grafos a processamento de sinal, aprendizagem profunda e aprendizagem quântica. Efetuam-se comparações e propostas de melhoria aos algoritmos existentes, e desenvolvida uma nova abordagem para abordar o tema de deteção de falhas. Com a retrospeção das ferramentas contemporâneas de maior sucesso, e pela oferta de uma nova abordagem, o objetivo deste trabalho é fornecer aos operadores de instalações fotovoltaicas o aumento na fiabilidade e eficiência dos seus sistemas. Além disso, há a possibilidade de que a ferramenta desenvolvida seja aplicável para outros problemas de coesão de dados, impactando positivamente os diversos tipos de domínios de sistemas orientados a dados.



# Abstract

The increase in utility-scale photovoltaic power plants has led to the need for effective methods for detecting and classifying component faults, which can have significant economic impacts. This work assesses the current state of fault detection and state estimation tools in the field of PV systems, focusing on understanding how these tools work and identifying their strengths and limitations. It is seen that machine learning makes up the majority of state-of-the-art fault detection and classification algorithms. Still, many fields have contributed to this problem, from graph theory to signal processing, deep learning, and quantum machine learning. Consequently, this work compares and proposes improvements to existing approaches or a novel technique developed to address this issue. By examining the most successful tools to date and offering new solutions, the intention is to help PV plant operators improve the reliability and efficiency of their systems. The developed methodology is also expected to become a generalistic data cohesion algorithm, positively impacting other data-driven fields.



# Agradecimentos

Aliquam id dui. Nulla facilisi. Nullam ligula nunc, viverra a, iaculis at, faucibus quis, sapien. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Curabitur magna ligula, ornare luctus, aliquam non, aliquet at, tortor. Donec iaculis nulla sed eros. Sed felis. Nam lobortis libero. Pellentesque odio. Suspendisse potenti. Morbi imperdiet rhoncus magna. Morbi vestibulum interdum turpis. Pellentesque varius. Morbi nulla urna, euismod in, molestie ac, placerat in, orci.

Ut convallis. Suspendisse luctus pharetra sem. Sed sit amet mi in diam luctus suscipit. Nulla facilisi. Integer commodo, turpis et semper auctor, nisl ligula vestibulum erat, sed tempor lacus nibh at turpis. Quisque vestibulum pulvinar justo. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Nam sed tellus vel tortor hendrerit pulvinar. Phasellus eleifend, augue at mattis tincidunt, lorem lorem sodales arcu, id volutpat risus est id neque. Phasellus egestas ante. Nam porttitor justo sit amet urna. Suspendisse ligula nunc, mollis ac, elementum non, venenatis ut, mauris. Mauris augue risus, tempus scelerisque, rutrum quis, hendrerit at, nunc. Nulla posuere porta orci. Nulla dui.

Fusce gravida placerat sem. Aenean ipsum diam, pharetra vitae, ornare et, semper sit amet, nibh. Nam id tellus. Etiam ultrices. Praesent gravida. Aliquam nec sapien. Morbi sagittis vulputate dolor. Donec sapien lorem, laoreet egestas, pellentesque euismod, porta at, sapien. Integer vitae lacus id dui convallis blandit. Mauris non sem. Integer in velit eget lorem scelerisque vehicula. Etiam tincidunt turpis ac nunc. Pellentesque a justo. Mauris faucibus quam id eros. Cras pharetra. Fusce rutrum vulputate lorem. Cras pretium magna in nisl. Integer ornare dui non pede.

David Freire



*“You should be glad that bridge fell down.  
I was planning to build thirteen more to that same design”*

Isambard Kingdom Brunel



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Fault detection in Utility Scale Photovoltaic Plants</b>	<b>5</b>
2.1	Utility-Scale Photovoltaic System's Architecture . . . . .	5
2.2	Faults in Photovoltaic Systems . . . . .	7
2.3	Modeling photovoltaic's physical/electrical behavior . . . . .	10
2.3.1	The five-parameter model . . . . .	11
2.4	Data Analysis and Characteristics of Photovoltaic Faults . . . . .	12
2.5	Literature on Fault Detection and Classification for Photovoltaic Systems . . . . .	12
2.5.1	Statistical and Signal Processing Algorithms . . . . .	13
2.5.2	Machine Learning Algorithms . . . . .	14
2.5.3	Deep Learning Algorithms . . . . .	16
2.5.4	Proposed method's scope . . . . .	16
<b>3</b>	<b>Preliminary Work Plan</b>	<b>19</b>
3.1	Academic and Industrial Setting . . . . .	19
3.2	Work timeline . . . . .	19
3.3	Development strategies . . . . .	20
3.4	Final expectations . . . . .	21
<b>4</b>	<b>CellTAN: Cellular Time Activation Networks</b>	<b>23</b>
4.1	Glossary . . . . .	25
4.2	The Cell . . . . .	25
4.2.1	Principles . . . . .	25
4.2.2	Processes and Data . . . . .	26
4.2.3	Inputs and Outputs . . . . .	28
4.2.4	Computing Activations . . . . .	30
4.2.5	Trust . . . . .	32
4.2.6	Connections . . . . .	39
4.2.7	Unconformities and State . . . . .	40
4.2.8	Application Plugins . . . . .	41
4.3	The Hub . . . . .	42
4.4	Expected Network Behavior . . . . .	42
4.5	Implementation . . . . .	43
4.5.1	Code and Infraestructure . . . . .	43
4.5.2	Communication Protocol . . . . .	43
4.5.3	Cell configuration and deployment . . . . .	44

<b>5 CellTAN Application</b>	<b>45</b>
5.1 Case study . . . . .	45
5.1.1 Data analysis . . . . .	46
5.1.2 Data Cleaning . . . . .	49
5.1.3 Photovoltaic Plugin . . . . .	50
5.1.4 CellTAN Configuration . . . . .	50
5.1.5 Simulation and Results . . . . .	51
5.1.6 Scaling up . . . . .	51
<b>6 Conclusion and Future Work</b>	<b>53</b>
6.1 Addressing the research questions . . . . .	53
6.2 Objectives reached . . . . .	53
6.3 Potential applications . . . . .	53
6.4 Future work . . . . .	53
<b>A CellTAN development</b>	<b>55</b>
A.1 Statistical tests for measure of association . . . . .	55
A.1.1 Pearson's chi squared test . . . . .	55
A.1.2 Fischer's exact test . . . . .	55
A.1.3 Odds ratio . . . . .	55
A.1.4 Phi coefficient . . . . .	56
A.1.5 Contingency coefficient C . . . . .	56
A.1.6 Theil's U . . . . .	56
A.2 Technology stack . . . . .	57
A.3 Cell configuration . . . . .	57

# List of Figures

2.1	Representation of utility-scale PV plant components and some possible faults. . . . .	6
2.2	Typical data flow of utility-scale PV power plants. . . . .	7
2.3	"Failures in grid-connected PV systems." . . . . .	9
2.4	"Circle chart related to the module defects in the 5 plants (over the total number of failures)." . . . . .	10
2.5	Single-diode model for photovoltaic modules. . . . .	11
2.6	Representation of some of the methodologies employed in fault detection for PV systems. . . . .	13
3.1	Dissertation work timeline represented by a Gantt chart. . . . .	20
4.1	Simples CellTAN Network of two cells cooperating. . . . .	23
4.2	The cell's core sequence of processes (colored orange) and flow of data (colored blue for private and green for public attributes). . . . .	26
4.3	Classical set, triangular fuzzy number, Gaussian distributions, represented as membership functions. . . . .	28
4.4	Visualization of the effect of time decay on $x_{lower}$ and $x_{upper}$ . . . . .	30
4.5	Visualization of a self-similarity extraction example. . . . .	32
4.6	Trust measurement evolution for different methods with $a = 30, b = 10, c = 10, d \in [0, 100]$ . . . . .	35
4.7	Trust measurement evolution for different methods with $a = 30, b = 10, c = 10, d \in [0, 3000]$ . . . . .	36
4.8	Trust measurement evolution for different methods with $a = 30, c = 10, d = 10, b \in [0, 1000]$ . . . . .	37
4.9	Trust measurement evolution for different methods with $a = 30, b = 10, d = 10, c \in [0, 1000]$ . . . . .	38
4.10	Trust measurement evolution for different methods with $b = 30, c = 10, d = 10, a \in [0, 1000]$ . . . . .	39
4.11	Algorithm for deciding nonconformity severity (intrinsic and extrinsic). . . . .	41
5.1	Inverter AC side power from 2020 to 2022, used for the knowledge base. . . . .	46
5.2	Inverter AC side power from 2023-01-01 to 2023-01-05, used for testing. . . . .	47
5.3	Pair plot of AC power from both inverters (2020 to 2022), using scatter (left) and KDE (Kernel Density Estimation) (right). . . . .	47
5.4	Pair plot of AC power from both inverters (2023), using scatter (left) and KDE (Kernel Density Estimation) (right). . . . .	48
5.5	... . . . .	48
5.6	... . . . .	48
5.7	... . . . .	49

5.8	...	50
A.1	Technology stack of the Cell and Hub of CellTAN.	57

# List of Tables

2.1	Comparison of literature that inspired this work. . . . .	17
4.1	Contingency table representing the activation intersection of two cells. . . . .	33
5.1	Available variables from two inverters and a satellite. . . . .	45



# Abreviaturas e Símbolos

ANN	Artificial Neural Network
CNN	Convolutional Neural Network
CXN	Cell Complex Neural Network
DC	Direct current
DL	Deep Learning
DNN	Deep Neural Network
LSTM	Long short-term memory
MCD	Minimum Covariance Determinant
ML	Machine Learning
PV	Photovoltaic
RBFNN	Radial basis function neural network
RMM	Recurrent Neural Network
SC	Short Circuit
SRC	Sparse Representation Classifier
STC	Standard Test Conditions
SVM	Support Vector Machine



# Chapter 1

## Introduction

The XIX century marked a significant shift in the world's perception of energy resources as the desire to invest in renewable energy sources to power modern societies grew. This transition was driven by the need to reduce dependency on fossil fuels, mitigate the effects of global warming, and slow climate change. Renewable energy sources offer a range of benefits, including reduced greenhouse gas emissions, increased energy security, and air quality. Solar photovoltaic energy is a desirable renewable energy source due to its abundance, accessibility, and environmental benefits. While solar photovoltaic energy has proven to be both cost-efficient and environmentally friendly, it also comes with unprecedented challenges, such as its intermittent nature, low electrical inertia, complex forecasting, and geographic-dependent operating conditions. Despite these challenges, recent reports [1] show that the economic benefits of investing in renewable energy outweigh the complications, as there is an increasing global investment trend in these sources.

The general construction of PV farms, particularly on the utility-scale, has led to a need for effective maintenance and monitoring to ensure maximum efficiency and operational reliability. Towards this, various algorithms and routines are used to monitor the state of PV farms and identify any potential issues that may arise. Fault detection is crucial to this process, allowing PV farm operators to identify and address problems quickly. Detecting faults and identifying the necessary steps can prevent or minimize downtime and ensure optimal performance. Given the importance of maintaining high levels of operation, knowing if action is needed to restore or fix components from an anomalous scenario is desirable for reducing investment risk and maximizing profits.

Integrating intermittent energy resources into modern electric grids has led to stricter requirements for connecting such power systems to ensure safe grid operating conditions. As a result, companies that own or plan to build photovoltaic farms must comply with these requirements and have adequate power electronics and monitoring/control capabilities. Failure to meet these requirements can result in sanctions or fines for the responsible party, as well as potential impacts on system availability, asset value, and disturbance propagation to the grid. To minimize these risks and maximize the value of their assets, companies may opt to implement fault detection and state estimation tools. These tools allow for the early detection and resolution of potential issues and can prevent or minimize downtime. The need to create or improve existing fault detection and

state estimation tools, and the search for the most effective methodologies for addressing these issues, drive research in this field.

Having laid the basis for why there must be system behavior assessment in utility-scale PV plants, it is necessary to understand what business concepts are crucial to this field. In the course of this work, the presented topics will go over the following questions:

- What components mostly fail in photovoltaic power systems?
- What is the average frequency of faults?
- What fault detection/state estimation tools exist for photovoltaic power systems?
- What are the most successful ones?
- What's their structure? Are they mostly centralized or decentralized?
- What are their computational costs/efficiency?
- What is the expected magnitude of precision and confidence?
- Which key performance indicators can evaluate the success of these tools?
- What are their implementation difficulties?

With these questions uncovered, the main objective is to adapt or design a novel algorithm/approach to fault detection based on modern artificial intelligence solutions. However, this can be split into finer goals:

- Identify and study existing fault detection tools for photovoltaic power systems.
- Adapt or develop a new tool.
- Apply and test the new tool in real case study PV assets.
- Validate the developed methodologies by comparison to reference tools.

Before reviewing state-of-the-art fault detection tools, types of failures in photovoltaic systems need to be understood: find which components usually fail, which ones fail more often, and how often. For this, it is necessary to understand such components' physical and electrical properties and the modeling techniques used to characterize them. There will be an assessment of utility-scale power plants architecture through literature, alongside the detection objective of state-of-the-art fault detection tools applied in this field. Then, there shall be an extensive analysis and review of what tools have been designed and used in this field. In this step, critical evaluation of the literature is a must for understanding the tool's scope, ease of implementation, and understanding that the data sets available for this work are compatible. Having selected the most prominent ones, they're to be qualitatively and quantitatively compared to each other in their application context so that

the results allow objective evaluations. This process requires implementing these tools, following the guidelines in the respective article/book/report, verifying their metrics, and checking if the achieved results resemble the same as the literature suggests. It will require gathering data sets, which can either be artificially generated through simulation or provided by an enterprise that services photovoltaic plant owners.

There's a desire that, in the end, the developed work helps achieve an improved method for fault detection and state estimation in photovoltaic power systems, resulting in a production-ready software application agile enough to deploy for multiple PV assets. It's intended that the algorithm specializes in data cohesion as a means of anomaly inference, allowing asynchronous and self-healing data transfers between the considered components. Depending on the new algorithm's characteristics, it could result in an approach capable of generalization and application to other engineering systems, benefiting more than just PV systems. No matter the chosen methodology, fault detection will, in most cases, result in an economic benefit, catastrophe prevention, and safety increase.



## Chapter 2

# Fault detection in Utility Scale Photovoltaic Plants

### 2.1 Utility-Scale Photovoltaic System's Architecture

Utility-scale photovoltaic (PV) power plants are large-scale systems connected to the electrical grid, having installed capacities ranging from kilowatts peak (kWp) to megawatts peak (MWp). These systems typically consist of many PV panels interconnected through power electronics to aggregate and inject power into the grid. The number and type of components in a PV power plant depend on the plant's scale and topology, with different configurations possible for large-scale applications, including central inverters, string inverters, and multi-string inverters [2]. The physical installation of PV modules can include solar tracking apparatuses, such as single and dual-axis trackers [3], which add to system complexity and change production behavior. Understanding the architecture and components of PV power plants is vital for designing, operating, and maintaining these systems, as it helps optimize their performance and reliability.

Figure 2.1 presents a typical utility-scale PV plant architecture using the central inverter (or possibly multi-string inverter) configuration. It is noticeable that many system components may fail in one or more ways, which is why monitoring and fault detection algorithms are essential to maintain state estimation. The main subsystems considered in this work are the following:

- Solar photovoltaic panels (with or without bypass diodes).
- Tracking mount.
- Electrical cabling.
- Inverter(s) (mostly with Max Power Point Trackers).
- AC Transformer(s).
- Protection components (circuit breakers, fuses, surge protectors, etc.)

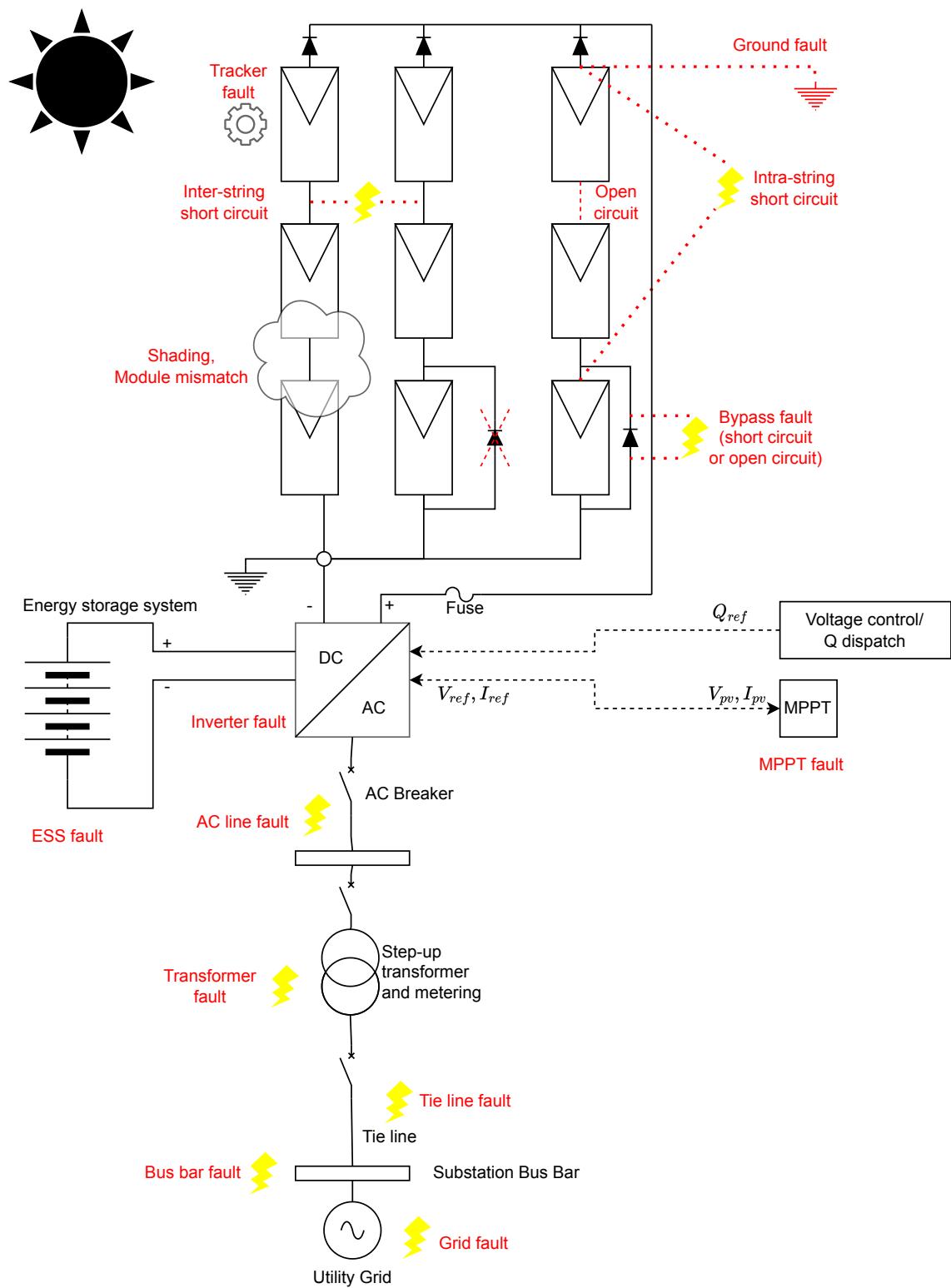


Figure 2.1: Representation of utility-scale PV plant components and some possible faults.

Most of these components have intrinsic variables, such as voltage and current values, that can help determine their operation states. Given that the utility grids (and the associated electricity

market) integrate large-scale PV assets, some of the before-mentioned components require constant monitoring and control, achieved with adequate embedded systems and sensor infrastructure [4]. Since monitoring utility-scale PV assets relies on the investment and technologies employed, engineers must consider data availability when developing data-driven algorithms. Thanks to the continuous advancements in communication technologies, namely in IoT (Internet Of Things), data acquisition is becoming faster, more reliable, and more precise. Not only is this fundamental for real-time asset assessment, but it also allows better training of fault detection algorithms. However, on the industrial scale (in the order of MWp production), having sensors embedded in every PV module comes with a high economic cost. Inverters are the components that usually possess monitoring capabilities, though the grid-tie connection should also be equipped with sensors. These can be considered the primary sources of information from utility-scale PV plants, with the most accurate, fast, and reliable data acquisition.

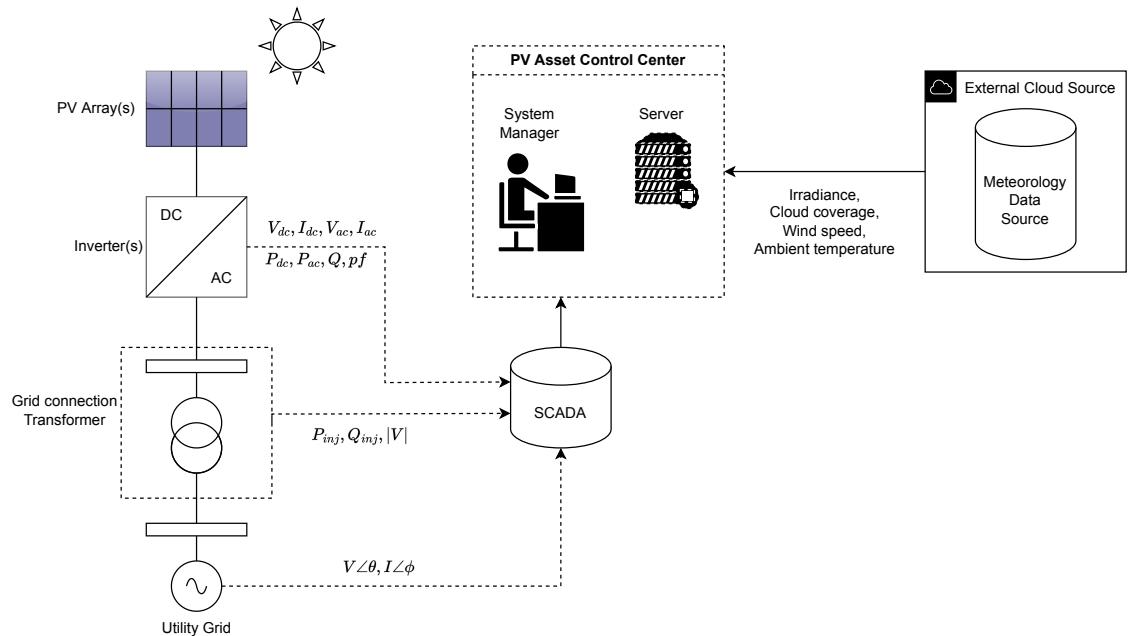


Figure 2.2: Typical data flow of utility-scale PV power plants.

Figure 2.2 represents a simplified data flow representation of a grid-tied PV system's most commonly available state variables, with most of them suggested by the IEC 61724 standard [5]. An external meteorological data source is defined since the PV system manager usually needs climate information for (at least) forecasting purposes.

## 2.2 Faults in Photovoltaic Systems

Several types of faults can occur in utility-scale photovoltaic (PV) power plants, which impact the performance and reliability of the system negatively. Unfortunately, some are very challenging

to detect and protect the electrical installation against, thus requiring sophisticated detection algorithms [6]. Besides the economical price, their occurrence may even cause safety hazards, such as fires [7], thus the urgency in detecting or preventing such events early.

According to [6], these faults can fit into three categories: electrical, mechanical, and environmental. Electrical malfunctions include short circuits, open circuits, and inverter failure, affecting the PV panels' power output and the system's overall efficiency. Mechanical faults include broken panels, damaged cables, and defective inverters, which can lead to system downtime and reduced performance (although not mentioned, solar tracker failures could also belong in this category). Environmental faults include extreme weather events, such as hail or strong winds, which can damage the PV panels and other components [8].

The authors in [9] cover a comprehensive review of most types of faults studied in the ambit of detection and classification algorithms. However, authors in [10] have a more succinct fault categorization that better fits this work's scope. They categorize all the major PV system faults into either DC-side or AC-side. Figure 2.3 represents this detailed categorization with a tree-like structure.

Although also prone to failure, most literature on fault detection and classification for photovoltaic systems does not encompass solar tracking faults: most studies cover fixed PV systems. The supervision and assessment of these subsystems' correct functioning can be sensor-based [11] or image-based. Some authors developed fault detection methods for these apparatuses [12], using image processing on aerial photography to determine modules' slopes. This category of failures should be better supported when developing electrical data-driven algorithms since they can significantly affect the system's efficiency. Hence, this work will attempt to include said fault category in the proposed fault detection methodology.

Throughout the literature [13], some of the most noted faults in the context of fault detection are:

- Shading: partial coverage of a PV array or module, temporary or not. It might result in a Hot Spot fault.
- Soiling: dirt accumulation, blocking sunlight from reaching PV Cells. It might also result in a Hot Spot fault.
- Short circuit: either line-line or line-ground.
- Open circuit: connection breakage between modules.
- DC arc fault: electricity plasma arc formed on broken connections.

According to a 2017 survey conducted on five utility-scale PV plants in Italy [14], the authors observed failure rates from <1% to 3% in the majority of plants and 81.8% in the worst scenario. The high failure rate of the latter had a demonstrated cause that originated from manufacturing mistakes: snail trails. Besides this phenomenon, hot spot faults and bypass diode faults/disconnections were among the most common.

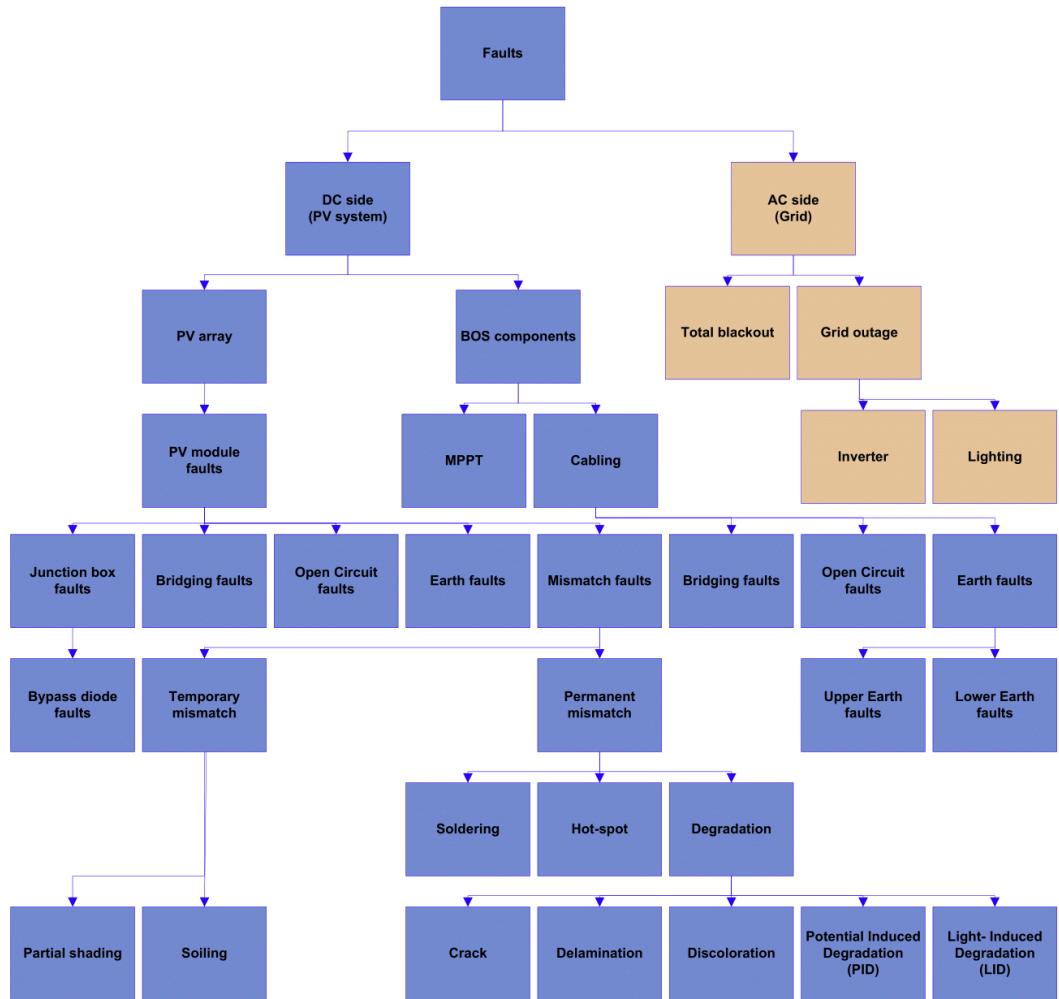


Image source and copyright: [6].

Figure 2.3: "Failures in grid-connected PV systems."

Alongside manufacturing failures, installation, planning, and other external effects can be the root cause for many of the presented faults [15].

Having the distribution of fault types from real-life scenarios is quite helpful for formulating fault detection algorithms. It allows for better generation/selection of training data and class decisions. In figure 2.4, it is possible to observe the failure type distribution for 24.254 inspected modules. Soiling, shading, and mechanically related failures were not as prominent, with only a group share of around 6%. It is relevant to note that discoloration represents almost a quarter of all faults.

Although the study had a limited geographic scope, with only a few power plants diagnosed, it allows for a more realistic view of the common scenarios encountered in typical operational ground-mounted utility-scale PV power plants.

Due to the difficulty of classifying some of these faults, given their similarity on the consequent effect in the system, it will be seen in further sections that most fault detection algorithms only

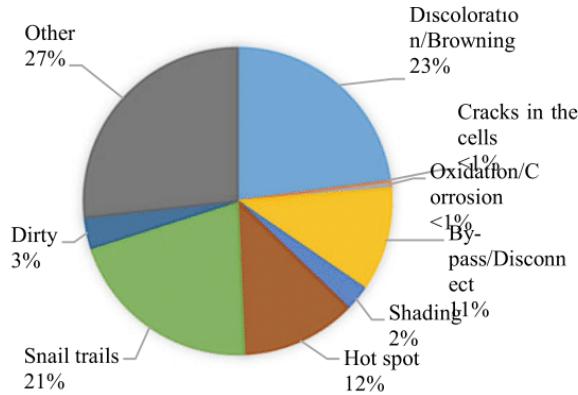


Image source and copyright: [14].

Figure 2.4: "Circle chart related to the module defects in the 5 plants (over the total number of failures)."

endeavor to classify between two to five types of reviewed faults.

## 2.3 Modeling photovoltaic's physical/electrical behavior

Photovoltaic cells are the fundamental components of photovoltaic panels. They are made from semiconductor materials like silicon and absorb photons that generate electric current. Their electrical behavior is characterizable using the current-voltage (I-V) equation 2.1. This equation, which represents a fundamental relationship governing the operation of PV cells, can be used to predict their performance under various operating conditions, such as differing solar irradiance and temperatures.

$$I = I_{ph} - I_d \times (e^{\frac{q \times (V_{pv} + I_{pv} \times R_s)}{n \times k \times T}} - 1) - \frac{V_{pv} + I_{pv} \times R_s}{R_p} \quad (2.1)$$

$I_{ph}$  (A) is the light-generated current;  $I_0$  (A) is the reverse saturation current;  $V_{pv}$  is the module's terminal voltage;  $I_{pv}$  is the module's output current;  $R_s$  ( $\Omega$ ) is the series resistance;  $R_p$  ( $\Omega$ ) is the shunt resistance;  $n$  (adimensional) is the diode ideality factor;  $k$  (J/K) is the Boltzman constant;  $T$  (K) is the cell temperature;  $q$  (C) is the electron charge;

For state estimation, it is crucial to accurately model PV modules' performance from the DC side of power converters. This information is vital for designing and optimizing PV power systems, as it enables predicting PV module performance under different conditions, as mentioned before. Accurate PV module models are also essential for state estimation and fault detection, as they provide critical information about their health and performance, allowing for early identification of potential issues. Moreover, they can be used to optimize the control and operation of PV power systems, improving their efficiency and reliability [13].

Physical and empirical models broadly categorize the several state-of-the-art methods for modeling photovoltaic modules [13]. Physical models lie on the fundamental physical principles governing PV modules' operation. They typically require detailed knowledge of the PV module's electrical and optical properties, such as its current-voltage (I-V) characteristics, spectral response, and temperature dependence. These models can accurately predict the PV module's performance under a wide range of operating conditions, but they may be complex and computationally intensive to implement [16]. On the other hand, empirical models are based on experimental data and are typically more straightforward to implement. However, they may not be as accurate as physical models, especially under conditions that differ significantly from those used to generate the experimental data (usually STC) [13]. Some examples of state-of-the-art physical models for PV modules include the single-diode model (the five-parameter model) and the two-diode model [17]. In contrast, one of the most used state-of-the-art empirical models is the Sandia model [13]. The choice of modeling method will depend on the specific application and the required level of accuracy and complexity; in some cases, there can be a combination of physical and empirical models.

In the case of utility-scale PV systems, detailed knowledge of the module's electrical and optical properties of empirical data may be limited, and building a model is only possible by recurring to the datasheet information. A complex model that requires more detailed information may not be feasible in such cases, and a simpler model that relies on fewer input parameters is more appropriate. Given the excellent trade-off between complexity and accuracy, the single-diode model suits this use case.

### 2.3.1 The five-parameter model

Figure 2.5 presents the single-diode model representation of the photovoltaic module. According to the five-parameter model, the unknown parameters are determined by fitting the model to experimental data or using data from the PV module's datasheet. The single-diode model can predict the PV module's performance under a wide range of operating conditions while maintaining reasonable accuracy. However, remembering that the single-diode model is a simplified representation of the PV module, it will have poor accuracy under certain situations compared to the more representative two-diode model [17].

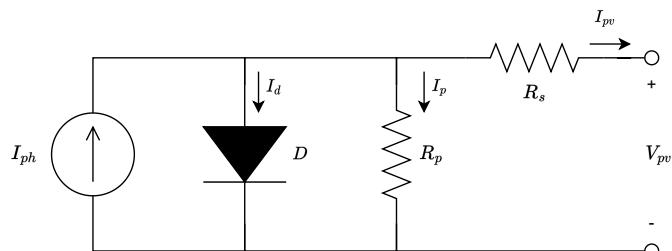


Figure 2.5: Single-diode model for photovoltaic modules.

## **2.4 Data Analysis and Characteristics of Photovoltaic Faults**

...

## **2.5 Literature on Fault Detection and Classification for Photovoltaic Systems**

The parent field of fault detection is anomaly detection (also known as outlier detection), a highly studied subject in the scope of statistics [18], applied in many scientific areas. Classification is also well-studied in this field, with applications in numerous scientific contexts, from medical diagnosis to airport safety [19]. Consequently, adaptations of generic tools and ad hoc methodologies have originated to aid in solving fault detection and classification problems in photovoltaics.

According to [4], the tools dedicated to PV fault detection and state estimation mostly come from mathematical/statistical methodologies, machine learning, and deep learning applications. Regarding the three general problem-solving principles mentioned before, it's known that machine learning and deep learning are the most popular and successful ones for recent applications that ought to solve complex problems. However, this categorization is somewhat limited, with contemporary literature suggesting an abundance of developed methodologies from different backgrounds, thoroughly reviewed in [9] and [10]. In [9], the authors consider two principal fault detection and classification algorithm branches: image-based and electrical-based; while [10] also distinguishes numerical-based techniques. Image-based refers to aerial or visual capture of the PV array by photography and thermal imaging, commonly used along with artificial intelligence algorithms for assessing the photovoltaic module's state. Although the contribution and importance of such methods are appreciable, this work will mainly focus on the electrical-based and numerical-based ones, as the use case of the developed tool is bound to this type of data.

Categorizing methodologies becomes fuzzy, considering that some literature mixes physical behavior models with machine learning, statistics, and signal processing. Figure 2.6 is an attempt to present a structure inspired by the review made by [9], [10], and this work's, with a focus on the more relevant techniques (for this work's scope). Hybrid models are ubiquitous since combining robust statistical, signal processing, ML, or DL models and PV's electrical characterization can achieve more remarkable results. Hence, a better representation than figure 2.6 would be an incomprehensible mesh of connections representing the permutations between category aggregation.

To not wander in the literature, there must be a decision on which methodologies to revise. The developed tool in this work must meet certain real-life constraints, such as data availability, frequency, accuracy, PV system configuration, and context. Therefore, the (qualitative) potential estimation for each methodology will be based on the capability of adapting the proposed algorithms to the same expected restrictions. This evaluation process confines the methodology review to emphasize the ones thought to be most capable of implementation in a real scenario. Therefore, the following sections will not cover an extensive literature review, as it is not intended to repeat

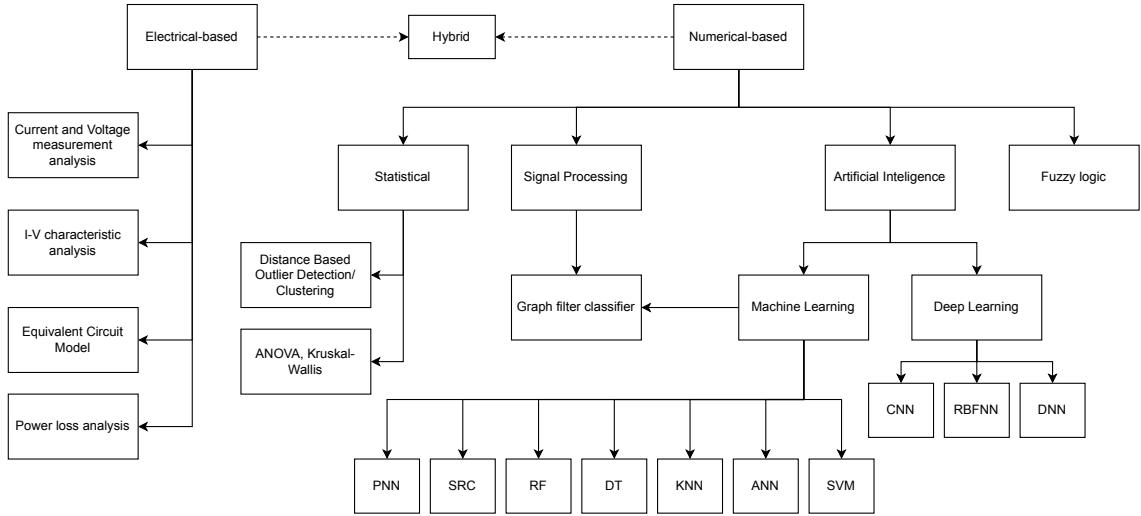


Figure 2.6: Representation of some of the methodologies employed in fault detection for PV systems.

the works of [9] and [10], only presenting interesting or adequate methodologies related to this work's scope.

### 2.5.1 Statistical and Signal Processing Algorithms

Statistical methodologies look into historical data to find the characteristics of how samples relate to the population (interpolation). These methodologies yield good results in case studies of PV farms that have been logging data for a considerable time, suffering in the cases that do not. Therefore, they are limited in that it is required to have curated data sets of historical significance for relevant features of the studied systems.

The literature on statistical and signal processing fault detection algorithms for PV is mostly quite dated ([20], [21], [22]), given that more recent machine learning methods have become increasingly attractive in this matter. Nonetheless, anomaly (or outlier) detection statistical algorithms can be used for fault detection in PV systems by identifying unusual patterns or deviations from normal behavior in the data collected from the PV system. Distance-based methods, such as the Euclidean, Mahalanobis, and MCD-based distances [13], may be adequate. Although simple, these techniques might only work for detecting outliers in the context of PV systems if they are scale-invariant (due to the different magnitude in the system's state variables) and resilient to outlier contamination (which only MCD-based distance is capable of). In [22], the authors applied Analysis of Variance (ANOVA) and Kruskal-Wallis test for inverter failure detection, with the only downside of only being able to identify outliers in a sub-array resolution, i.e., not for specific string or module failures.

Some algorithms consider incoming data from PV systems as signals, allowing the adaptation of signal processing theory to develop ad hoc algorithms. Coming up with a relatively simple algorithm, the authors in [23] propose a power-based fault detection method that only requires delayed

samples of the PV array's power output and a threshold. It reasons that since the power output of PV systems can't vary beyond a given point, considering a very short-term period (milliseconds), significant perturbations in this variable can be associated with faults. Although the simplicity and ease of implementation, it's clear that the success of this method requires feeding the algorithm with relatively high-frequency data, which would only be feasible on-site (and with specialized monitoring equipment).

In [24], the authors successfully formulated a graph signal processing algorithm for fault classification that yields increasingly better results when there is a considerable amount of labeled data, although its training is only semi-supervised. The results outperformed other standard machine learning methods for the same training data, given 30% or more of labeled data. On another note, the data utilized came from the PVWatts [25] dataset, and the PV system is on a small scale (ASU testing facility [26]) possessing a monitoring density and capability that can be considered unrealistic for utility-scale. This same data source is present in many other reviewed works.

The authors in [27] displayed another excellent use for graph theory, although not specifically for fault detection: they implemented a consensus-based distributed approach to minimize the impact of noise in acquired data from the PV array. By formulating a data propagation algorithm that resulted in measurement convergence, they achieved higher accuracy for state estimation.

With both graph theory-based algorithm proposals, this field sparks interest in its usage for the upcoming formulated methodology, given that it would be desirable to achieve an algorithm that features fault detection alongside data consensus.

### 2.5.2 Machine Learning Algorithms

Machine learning is the trending way of solving increasingly complex and non-linear problems, as neural networks (or other learning structures) can better model complex, non-trivial, and nonlinear relations between data. Still, they are as good as the training data, with many designs requiring a lot of representative learning examples to achieve good results. Their output can also be very obfuscated (depending on the technique), meaning that many methods do not allow a direct interpretation of the relationship between inputs and outputs. This "black-box" characteristic, specifically of neural networks, is considered a disadvantage. Besides, extrapolating data remains a challenge when classically using these structures. Still, they have immense applications for PV systems, from MPP (Max Power Point) estimation to power forecasting, soiling, and fault prediction.

In [28], the authors utilize an ANN to classify short circuit and hot spot faults. This algorithm achieved an outstanding 98.4% classification accuracy, yet the data originated from *MatLab/Simulink* simulations and only considered two classes of faults. Because the inputs were the variation of voltage and current ( $\frac{dV}{dt}, \frac{dI}{dt}$ ), the algorithm required data sampling with relatively high frequency (>5Hz). The present work will not regard such methodologies as background for the upcoming tool since requiring high-frequency simulated data while covering only two fault types is quite far from a real utility-scale PV system scenario.

The trend of utilizing simulated data (sometimes without even added noise) has been a target of criticism in [29]. Accordingly, this work also emphasizes that the literature shows many proposed

ML (and other types of) techniques that fall into this concept, which makes selecting appropriate methodologies to base future work on a challenging task.

The proposed ANN solution in [30] is remarkable by the diversity of fault classification achieved: STC, short circuit, varying temperature, partial shading, complete shading, degraded modules, ground fault, and arc fault. It presents one of the most fault class coverage with high accuracy, considering the literature that utilizes synthetic noiseless data. Hence, the cyber-physical conceptualization and data preprocessing (clustering) demonstrated can be admired, but not forgetting that validation data came from a relatively unrealistic setting.

In [31], there is a captivating proposal of utilizing an autoencoder and pruned neural network to separate the tasks of detecting and classifying faults, which resulted in one of the most performant ML approaches in the literature. The algorithm classifies five states: degraded, shaded, soiled, short circuit, and STC, utilizing nine inputs representing voltage, current, power, and irradiance available from the MPPT, datasheet, or meteorological sources. While the neural network pruning adds complexity, it resulted in a better generalized and lighter-weight trained model suitable for faster detection times. Even though using data from a small-scale PV system, the presented algorithm and its assumptions may make it possible to adapt and implement in an industrial scenario.

Regarding performance, the work in [32] proposes a sparse representation classifier (SRC) that evaluates if the system has line-to-line or line-to-ground faults for varying operating conditions. Although a drop in accuracy occurred for extreme circumstances, it is impressive that the algorithm identifies faults in such varied operating conditions: 10 to 50 degrees ambient temperature, 200 to 1000  $W/m^2$  irradiance, 10 to 60 % of mismatch, and 0 to 25  $\Omega$  of fault resistance. The feature extraction step was also very impressive, which could be a determining factor in the method's performance. Unfortunately, this work also does not validate results with experimental data and only uses simulation as a source. However, the demonstrated computational performance, both in terms of training cost and utilization speed, its usage without the need for training for parameter tuning, the straightforward implementation, and consistent convergence, suggests the potential for this alternative in the face of other ML methodologies. The authors also emphasize that sparse representation might be utilized alongside different learning algorithms for classification, opening the door to many possible future implementations.

An exciting yet far-fetched proposal was made in [33], where a quantum neural network (QNN) is formulated for PV fault classification. The QNN was trained for predicting just two scenarios: faulty or standard, but required up to four days of training, resulting in 93.89% accuracy. For comparison, the classical ANN took twenty seconds to train and achieved 95.39% accuracy. Although the methodology showcases the potential of quantum computing for this field, its preliminary results still distance itself from the traditional methods.

An abundance of ML methods have been tested and reviewed in this field ([9],[10]), utilizing structures such as SVM, KNN, RF, etc. Nonetheless, the results of [31]-[32] sparked the most interest in this work's scope.

### **2.5.3 Deep Learning Algorithms**

The field of deep learning is a branch of machine learning, with the term "deep" referring to amplified machine learning structures that ought to understand data patterns through more complex and intertwined artificial neuron connections. A simple example of a deep learning model would be the design of an artificial neural network with multiple hidden layers (DNN), with the intuition that each of these "extra" layers achieves feature/pattern recognition in a cascade. Other DL structures include the LSTM, CNN, and RBFNN. They have been explored alongside classical machine learning techniques for PV fault detection, although the known disadvantage is a usually high computational cost and relatively tricky implementation. These techniques are typically applied to image-based solutions [34] since they require classification based on 2D data from various image acquisition equipment [35], [36]. Given the 1D characteristic of raw electrical data, little literature considers these techniques for fault detection, as it implies an extra step of increasing dimensionality. However, there are some promising results in doing so [29].

In [29], not only is a DL technique presented for fault detection and classification, but there is also the best attempt at comparative evaluation against other methodologies. As mentioned, much of the literature presents results solely based on particular datasets comprising simulated noiseless data, invalidating any significant quantitative comparison.

Authors in [37] use a CNN model based on the pre-trained AlexNet for classification and feature extraction, allied with a classical ML model also for classification. The classified faults were arc fault, partial shading, open circuit, and short circuit. While the experiments utilized simulation data, adding noise and an abundance of heterogenous operating conditions better resembled a real scenario. Considering the same noisy data, other tested methodologies present 22-70% average accuracies, with the proposed fine-tuned AlexNet CNN reaching a maximum of 70.45%. This work presents one of the best benchmarks in the literature, with decent coverage of other state-of-the-art ML and DL algorithms, while demonstrating the most realistic results and a sophisticated methodology proposal.

### **2.5.4 Proposed method's scope**

While classical fault detection resides in the synchronous and direct evaluation of state and climate variables, realistic industrial scenarios can have data from various types, sources, and acquisition rates. It's also important to realize that monitoring equipment can register erroneous information, and current communication technology is also susceptible to delays and data loss. With this in mind, recent developments in the intelligent composition of deep learning structures aligned with graph theory spark some interest in their application to this field, such as the new deep learning technique named Cell Complex Neural Networks [38]. The motivation for choosing such a structure comes from its data propagation and consensus capability. The propagation techniques utilized in a CXN appeal to graph theory, dividing a system into other subsystems and components (nodes, also called cells in [38]) that share information. Even if the direct application of this structure might not be feasible or grant better results in the context of fault detection, its modification

to meet the scope's needs could result in a robust and efficient solution. Further investigation of this state-of-the-art tool will unroll throughout the development of this work in an attempt to adapt this knowledge to the PV fault detection field.

According to the reviewed methodologies, the proposed tool should pertain to the DL or the hybrid category since, while having a central component of DL, it may also require modeling the PV system's components. The intention of proposing such a novel approach is to contribute to the deep learning methodology ecosystem, explicitly formulated for electrical-based PV fault detection and classification. As mentioned, it aims at an asynchronous and online application, which differs from most current methods and presents a novel DL paradigm considering current knowledge. This work also desires to bring a comprehensive benchmark between popular methods (likewise [29]), utilizing a richer dataset with samples from tangible utility-scale PV assets, allowing accuracy assessment in a realistic scenario.

Reference and year	Data Source	Inputs	Proposed methodology	Classified Faults (alongside STC)	Validation data realism	Computational cost	Notes	Drawbacks
[29] 2020	Simulated PV System, added noise	Irradiance, Temperature, Short circuit current, Open circuit voltage, PV current, MPP current, MPP voltage, MPP power, Boost converter Maximum current, Voltage and power.	Pre-trained CNN (AlexNet) for feature extraction and classification	Arc fault, Partial shading, Fault during partial shading, Open circuit, Line to line SC	Moderate	High	Resilient against noisy data. Outperforms classical ML methodologies.	Requires data samples from the MPPT boost converter.
[32] 2020	Simulated PV System, no added noise	MPP voltage, MPP current, Short circuit current, Open circuit voltage, Irradiance	Sparse representation classifier	Line to line SC Line to ground SC	Low	Low	Very fast learning speed compared to classical ML structures. Straightforward implementation. Good feature extraction process.	Validation data was very idealistic. Only classifies line to line and line to ground faults.
[24] 2020	PVWatts dataset	MPP voltage, MPP current, Short circuit current, Open circuit voltage, Irradiance, Fill factor, Temperature, Gamma ratio, Maximum power	Graph signal processing	Shading Degraded modules Soiling Short circuit	High	Low	Semi-supervised, allows usage of unlabeled data for training. Better accuracy relatively to other ML methods for less labeled data. Low training cost.	-
[31] 2021			Autoencoder for detection and pruned neural network for classification			Medium	Separate the task of detection from classification, allowing for other combinations. Good performance method considering the algorithms complexity. Pruning creates an ANN less prone to overfitting.	Requires more complex training phase, for two different networks, and utilizing a dropout algorithm for pruning.

Table 2.1: Comparison of literature that inspired this work.

Table 2.1 represents a summary that compares four of the most inspiring reviewed proposals.



# Chapter 3

## Preliminary Work Plan

### 3.1 Academic and Industrial Setting

The present work unrolls at Faculdade de Engenharia da Universidade do Porto by a student employee of Enlitia [39]. One of the employer's requirements is to formulate the final algorithm in the Python programming language [40]. Therefore, we consider academic and industrial standards during the development of the tool.

Having the possibility of working with a company that provides artificial intelligence solutions for energy systems, there is ample availability of PV asset data from various clients, mainly from the Iberian Peninsula and other European countries. There will be a need to gather information from assets with historical significance and with the presence of faults. Although Enlitia leases PV asset data for this work, its origin will remain classified.

### 3.2 Work timeline

Figure 3.1 represents the expected chronological sequence of events. The most uncertainty comes from adapting the formulation of CXNs for PV fault detection and developing a Python application that implements its behavior, given that the know-how associated with this step is yet to be understood. There is also some doubt about the work time required to upscale the experimental environment toward actual deployment, knowing that the empirical setting of data science experiments may radically differ from software production environments. There will be an attempt to match the development environment with the implementation context, reducing the latter's risk.

According to the predefined workflow, the following chapters of this document will include:

- Deep analysis of Cell Complex Neural Networks and adaptation proposals towards PV fault detection and classification algorithms.
- The process of data gathering, defining data sources, and available information.
- Exploratory data analysis (EDA) and data mining (feature extraction, clustering, etc.).

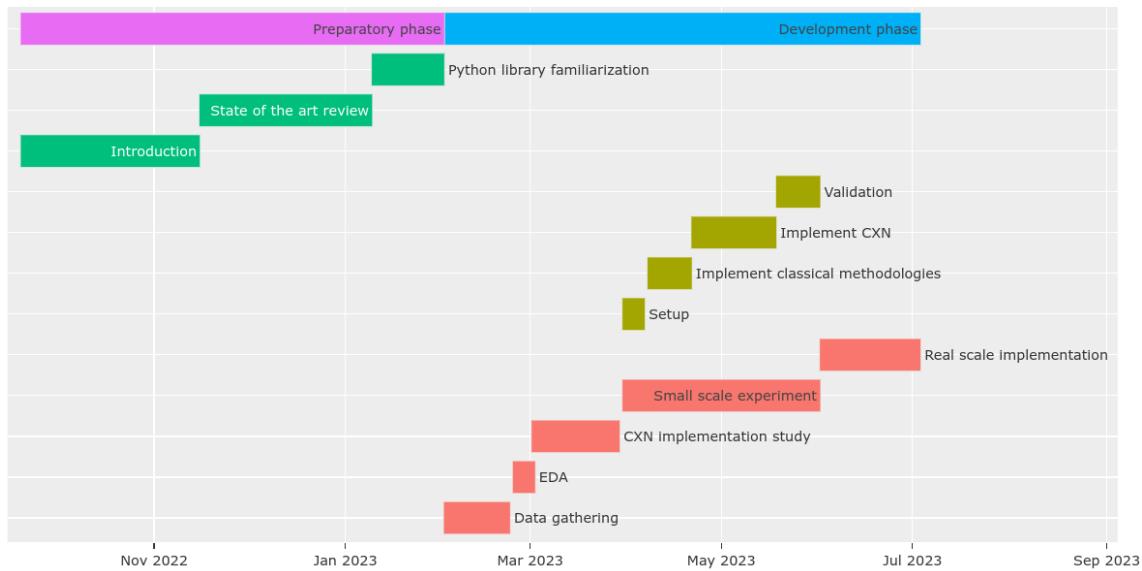


Figure 3.1: Dissertation work timeline represented by a Gantt chart.

- Setup context for the small-scale experiment: dataset construction, cleaning, normalization pipelines, etc.
- Implementation of classical ML methodologies.
- Implementation of novel CXN methodology.
- Result validation.
- Productization of the developed algorithm.
- Large scale deployment of the application.

### 3.3 Development strategies

Data gathering will be accomplished by combining as many data sources for the same variables as possible, such as power meters and SCADA measurements, local meteorological stations, and satellite data. Then, datasets will be organized into specific categories according to time resolution, acquisition method, availability, and accuracy. Geographical-based characterization might also be utilized to divide data even further. Data analysis, feature extraction, and clustering are essential in this step to harness valuable data from unlabeled samples.

To implement classical ML methodologies, as well as the proposed tool, Python libraries such as TensorFlow [41], Scikit-learn [42], and SciPy [43] are the finest selection. They facilitate the creation of neural networks and other learning structures, possessing many different implementations of scientific algorithms. Other utilities discovered during the development process will be referenced along.

Software engineering guidelines will rule the software architecture of the final application, which should follow its well-studied organization patterns. Even though it is not a focus of this work, paying attention to code implementation should result in a better application, scalability, and security.

### 3.4 Final expectations

The final application aims to implement online fault detection for operational PV assets. For performance reasons, its deployment architecture is a consideration to take in prior, such as its potential for parallelization. It should be as generalized as possible, as that would increase its correct functioning for different PV assets of various plant operators. In the end, it is expected to resemble the usability and robustness of a finished software product so that PV plant operators can reliably utilize its outputs. Its success will result in higher PV system efficiencies and safety, advancing the world's energy transition.



## Chapter 4

# CellTAN: Cellular Time Activation Networks

Distributed information systems characterized by time series data present various challenges, primarily due to their complex and dynamic nature. The sheer volume of data that must be processed and analyzed in real time is a significant challenge, leading to concerns over storage, computation, and scalability. Moreover, data quality issues such as missing or incomplete data and data heterogeneity arising from sourcing data from disparate sources with varying consistency and structure further exacerbate these challenges. Another critical challenge of these systems is handling the temporal aspects of time series data, requiring specialized pre-processing, feature extraction, and modeling techniques. The distributed nature of these systems further complexifies matters, with issues encompassing data synchronization and accuracy being a common concern. Furthermore, their implementation in real-world applications requires robust mechanisms for data security and privacy, which adds a layer of complexity to the design and implementation of these systems.

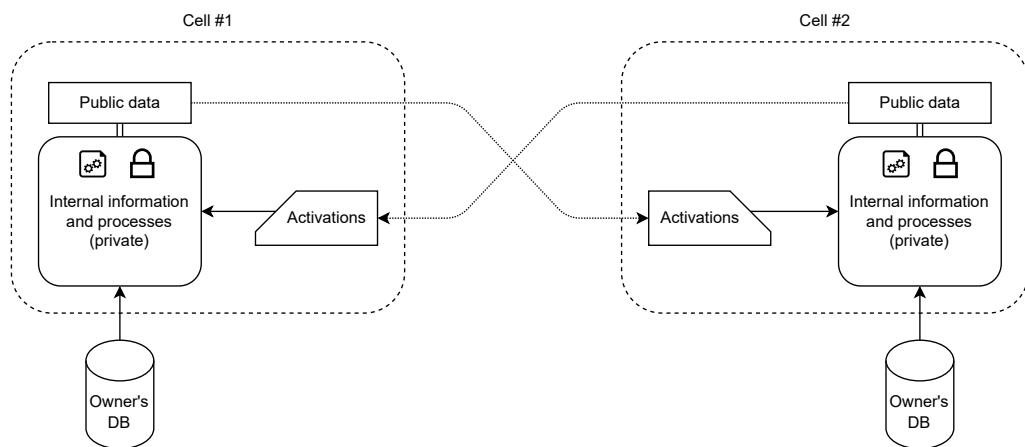


Figure 4.1: Simple CellTAN Network of two cells cooperating.

This chapter proposes a novel tool entitled CellTAN (Cellular Time Activation Network) that

undertakes those challenges. CellTAN represents sparse yet interconnected components that function independently, cooperatively, and asynchronously. Inspired by other effective mechanisms like GNNs, CXNs, and Weighted Cross-Connection Networks (WCCNs), CellTAN uses a graph-like structure to represent a network of components with nodes and connections. Following the introductory chapters, its primary purpose is to detect abnormal scenarios on PV systems. However, its generalized formulation introduces other valuable features which come naturally from fulfilling this goal. Such are state estimation, forecasting, and capturing the value in data from different PV asset owners without violating their privacy. For brevity's sake, we will unfold details about its benefits during the rest of this chapter.

Instead of tackling fault detection and classification in a classical centralized manner, which is already extensively showcased throughout the literature, this tool approaches this problem with a paradigm change: a distributed and asynchronous data coherence system. By having a virtual representation closely related to the physical form of sparse systems, we can leverage the relationships between components to assess their correct (or incorrect) operation. While initially designed for photovoltaic (PV) systems, the concepts of cells, connections, neighbors, time series, and uncertainty are universal and applicable to other fields such as biology, physics, and more. Thus, the potential for generic applicability sparks the interest to not bake specifics of PV systems directly into this tool, allowing its usage for other subjects. Figure 4.1 represents a minimal scenario for a network: only two cells. It showcases some terms specific to this tool that might be difficult to grasp initially, such as trust, events, and activations. Consequently, the glossary in 4.1 and detailed explanations throughout this chapter serve to clarify them.

Vast and scattered information across multiple agents is a common scenario faced by the industry of AI for energy systems, which cannot be aggregated due to privacy and confidentiality reasons. Nevertheless, its conjunction could have a lot of added value, given the similarity of certain assets: PV plants (as well as wind farms) from different owners in neighboring geographical regions. This information-sharing potential for AI algorithms motivates the development of a mechanism that communicates information between differently owned assets without any of the compromises above. However, as stated before, data acquisition in PV scenarios is scarcely synchronous and might only occur in equal time resolutions for some of the different components. The CellTAN addresses these issues with an instrument called **Activations**. It proposes a new way of communication that decouples from the needs of units, sampling rate, and synchronization, avoiding resampling, normalization, or even obfuscation (to protect privacy), by only transmitting time values. This mechanism is a core feature of the tool since it will be the means that will allow connecting different stakeholders' data, and section 4.2 develops this matter thoroughly. Likewise, succeeding sections formulate the working of the cell and its interactions within the network. Since it is the core component, understanding its behavior is crucial for a complete understanding of the tool.

## 4.1 Glossary

- **Knowledge base:** Refers to registered historical knowledge (samples) of time series variables.
- **Inputs:** Uniform fuzzy numbers (not necessarily, but it is the current choice) that represent one sample of the group of time series variables that define the state of a cell.
- **Outputs:** Similar to inputs, but obtained through some computations.
- **Time decay:** A process associated with increased uncertainty of variables over time.
- **Activations:** Timestamps of past occurrences on a knowledge base with a non-zero membership value against a set of recent inputs.
- **Trust:** A decimal number that represents how coherent two activations are with each other.
- **Hub:** The central component of the cell network, which facilitates its visualization, management, and expansion. It acts as the proxy agent between the cell's communication.

## 4.2 The Cell

### 4.2.1 Principles

A cell is an independent entity composed of **data** and **processes**. The idea is to abstract fundamental system components (e.g., inverters, MPPTs) into this virtual entity. With an added intelligence layer and featuring a few different processes, it assesses its current state based on all available internal and external information, adding value to the existing data acquisition and monitoring systems. As an individual part of the system, it follows a set of rules that define its intrinsic and extrinsic behavior. These rules address data privacy, request boundaries, and real-world operational limitations.

**Independence** During operation, independence on neighbors or other network entities for continuous processing of outputs results in a more robust system and increases cell availability. Thus, given any connection cutoffs, the cell shall be unbothered by its surroundings and continue operating in an isolated state. Isolation is not preferred, but enduring it until outside contact is re-established avoids shutdown and startup procedures.

**Selfish Computations** The cell is selfish in that it will not perform any computations based on the request of others. This aspect creates a fundamental layer of protection against overloading the infrastructure in which it is deployed, which also increases availability.

**Data accessibility** Although selfish in computations, the cell shall provide access to select data valuable to the network. However, not all internal data is shared, and public data shall not compromise the cell's privacy (more on this later).

#### 4.2.2 Processes and Data

The cell's core is essentially characterized by a main process loop that executes a sequence of actions periodically.

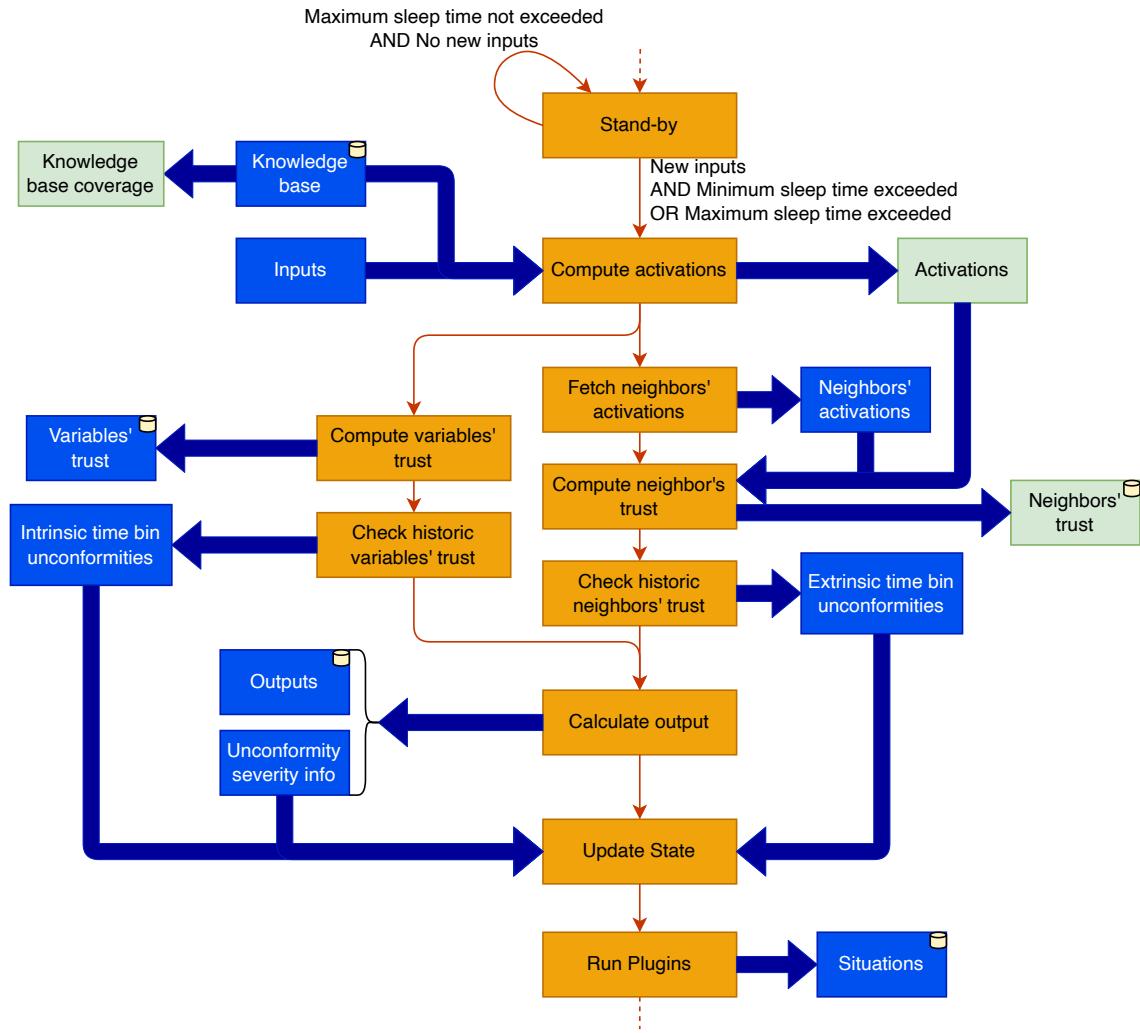


Figure 4.2: The cell's core sequence of processes (colored orange) and flow of data (colored blue for private and green for public attributes).

Figure 4.2 showcases the cell's core procedures and data attributes. The symbol on the top right of some attributes indicates that the cell should persist in its database. Regarding the color coding in figure 4.2, we can see that most of the cell's attributes are private. The network can access its activations, knowledge base coverage, and neighbors' trust. Conversely, the cell extracts neighbors' activations from the network and gives its owner exclusive access to inputs, outputs,

knowledge base, and others. The relationship between the owner and private cell data is bidirectional since it is also his responsibility to update the inputs (hopefully with an autonomous process). The cell's knowledge base may only represent an interface with a private database (thus a bidirectional connection), not the data itself, since the owner might prefer a centralized storage server instead of storing it within the hardware of the cell.

The cell starts in standby mode, which is only left when new information arrives or a predetermined amount of time has passed. This mode ensures no unnecessary computational burden of repeating all processes without new data. Nonetheless, defining a maximum sleep time ensures the cell's activations are periodically updated to reflect contemporary neighbors' trust and intrinsic uncertainty.

The process of computing activations relates to temporal similarity extraction, presented in section 4.2.4. It is associated with searching for similar behavior (of the inputs) in the cell's knowledge base, generating "time activations", which are the timestamps of similar occurrences (this is made publicly available).

Fetching neighbors' activations is an asynchronous process that uses the hub for communicating with other cells. This process only occurs after checking what neighboring cells have activations not older than an acceptable threshold (configured by the cell owner) and if their knowledge bases intersect with the cell's own. However, if a neighborhood of cells receives new inputs simultaneously, it would result in some trying to continue computations while others still do not have activations. Thus, this process checks for neighbors with updated activations up to three times, with a pause in-between, if the current amount does not reach a predefined threshold.

The cell executes a trust computation neighbor-wise and variable-wise. In practice, this procedure generates a decimal value between zero and one for each input variable and neighbor, which describes how much that variable is coherent with the rest and how much a neighbor is coherent with the cell (respectively). We nicknamed it "trust" to have a friendlier term used throughout the development and reasoning that unrolls in section 4.2.5. This indicator is stored in the owner's database so that it can perform an aggregated analysis, which relates to the process of "checking historical trust".

The "checking historical trust" procedure identifies unconformities in aggregated trust values during specific time windows and flags them when they fall below the owner's threshold. It records inconsistencies alongside information about the elements and time bins involved, raising the severity of unconformity accordingly.

For output computation, the cell uses intrinsic and neighbor activations to filter its knowledge base and find the new variables' domain (a process detailed in 4.2.3). This process also provides information about unconformity severity, since certain unwanted situations might occur when intersecting variables' and neighbors' activations (see section 4.2.7).

During the "updating state" stage, the cell checks if its internal state has changed since the previous loop. This state is closely related to the possible unconformities (4.2.7). It sends this information to the owner's database and the central hub if a change does occur.

In the end, the cell may run application-specific plugins. These are external additions that the owner might include for uncovering specific situations related to the cell's inputs. Further detailed in section 4.2.8, this mechanism binds system logic to the cell and takes advantage of its known physical properties and behavior.

### 4.2.3 Inputs and Outputs

Inputs and outputs are a view of the values that define a cell's variables at a given rolling timestamp. While inputs are directly associated with raw sampled data from the system (injected by the owner), outputs are a byproduct of internal processes. The latter should present more accurate information since it is based on internal and external data (ideally) and may be helpful for the cell's owner to assess its state.

We defined inputs and outputs using classical sets instead of crisp values, which will be crucial for the cell's inner workings. Representing the cell's variables in a fuzzy (or probabilistic) manner can better capture the inherent uncertainty in time series data. However, they are not limited to this representation, with fuzzy numbers or probability distributions as alternatives. Besides, they can be subject to a process called **time decay**, which ensures that the passage of time negatively affects uncertainty (more on section 4.2.3.1). This mechanism increases the robustness of the cell by acknowledging the value of time in assessing its current state.

Generalizing the concept of representing uncertainty, the following structures may characterize inputs and outputs:

- Classical set: simple uncertainty band (e.g., uncertainty up, down, relative, etc);
- Fuzzy number: generalized fuzzy number representation [44] (e.g., triangular fuzzy number  $(a,b,c;h)$ );
- Probability distribution: the distribution's characteristics (e.g., mean ( $\mu$ ) and standard deviation ( $\delta$ ) for Gaussian, the mean rate of occurrence ( $\lambda$ ) for Poisson, etc.);

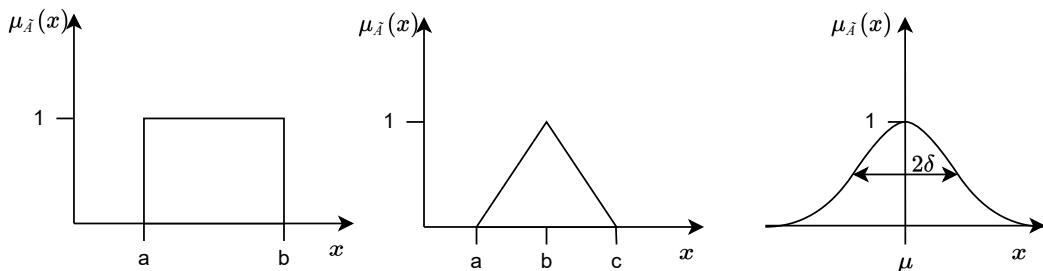


Figure 4.3: Classical set, triangular fuzzy number, Gaussian distributions, represented as membership functions.

Classical sets have pros and cons: many operations become more efficient than the alternatives, but we lose density information. Filtering historical data becomes trivial: a value lies within the

bounds of the interval (membership value of one) or does not (membership value of zero). This more accessible representation will benefit some cell processes, primarily in temporal similarity extraction (4.2.4).

#### 4.2.3.1 Time decay

In real-world dynamic systems that involve data acquisition, the certainty of the data collected tends to fade over time due to its dynamic nature. Typically, the most recent data is the most accurate representation of the system's current state, and as time elapses, the accuracy of previous data points decreases. As a result, it is crucial to consider the time dimension when analyzing dynamic systems and to develop methods that can account for the decay in data certainty over time.

As stated before, and towards considering the time dimension for the current state of a cell, we formulate a **time decay** method to ensure a more truthful and reliable assessment of the cell's current state. During the standby stage, this process ensures that the inputs and outputs suffer an increase in uncertainty, leading to an uncertain cell in a steady state (without new data input).

Consider the following example of converting a crisp value (5) and uncertainty ( $\pm 1$ ) to a classical set:

$$x = 5 \pm 1 \rightarrow [5 - 1, 5 + 1] = [4, 6]$$

To simulate the increase in uncertainty over time, we suggest the following equations (4.1 and 4.2), applied to each bound:

$$lower = lower - (lower - minimum) \times \frac{age}{decay} \quad (4.1)$$

$$upper = upper + (maximum - upper) \times \frac{age}{decay} \quad (4.2)$$

The *age* variable refers to the time difference between the present time and the instant of data acquisition. For implementation, the time unit considered is the 'second'. The *decay* is a parametrized constant (same unit as *age*) that defines the time it takes for the set to expand into its limits when starting from the median. It is chosen based on the characteristics of the variable, i.e., knowledge of its uncertainty over time. However, a good starting point is defining close to the data acquisition period so that the set expands entirely until a new value is acquired.

Figure 4.4 represents the expansion of the set throughout a period equal to the decay parameter (ten seconds). When the *x* variable reaches the age of ten seconds, its set representation transmits complete uncertainty: the bounds become similar to the variable limits ( $x_{min}$  and  $x_{max}$ ). We can see that the decreasing difference between the bounds and maximum/minimum values causes attenuation in the decay curve, as it displays a non-linear behavior. This behavior seems appropriate according to the reality of systems: as a variable becomes more uncertain with time, there is less potential for its uncertainty to increase.

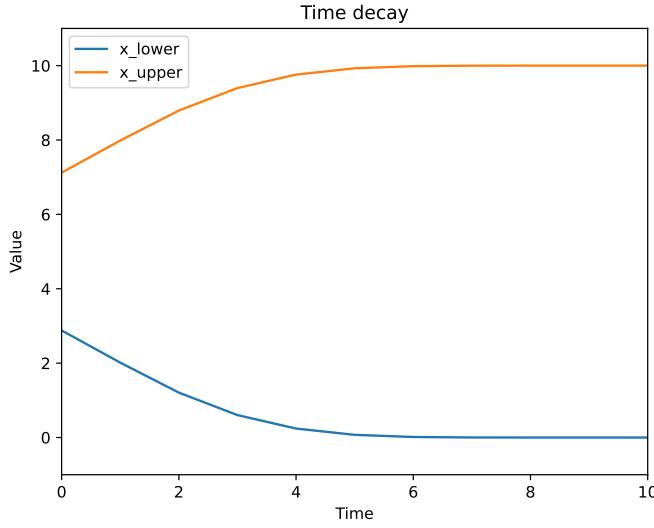


Figure 4.4: Visualization of the effect of time decay on  $x_{lower}$  and  $x_{upper}$ .

#### 4.2.4 Computing Activations

Earlier, we brought up the concept of temporal similarity extraction concerning computing activations without fully explaining it. This procedure involves identifying recurring patterns in a statistical base, i.e., finding past instances where we observe the current state to extract useful information about the system's behavior over time. By identifying historical periods with similar conditions, temporal similarity extraction can help assess the current state or predict future trends. This technique of historical search is not exclusive to this application, with demonstrations of use in other fields and for other purposes.

Using sets or probability distributions to filter historical data is one approach to simplify the process of temporal similarity extraction in multivariate time series data while also making it more robust against noisy or incomplete data. This approach assigns membership values to each data point in the time series based on the corresponding set (or distribution) generated from the current cell inputs. We form a group of similar past instances by eliminating samples past a certain threshold. This technique is one of the reasons for deciding to represent the inputs and outputs of the cell in a fuzzy manner.

The proposal for similarity extraction in the cell consists of taking the newly received values (inputs) of the cell's variables (coming from sensors or other data acquisition equipment, calculations, forecasting, etc.), generating a classical set, fuzzy number, or probability distribution from them, and then applying the bounds/membership function or probability density function to the knowledge base (see figure 4.5). If historical samples are associated with membership values, there may be a process for determining outputs by combining data and membership values.

The initial choice is to use classical sets since filtering history becomes trivial and efficient: restrict the knowledge base to samples where all variables belong to the corresponding interval. Generating outputs with these can be as simple as constructing a new set based on the bounds of

filtered knowledge (samples with membership equal to one). When filtering historical data with two or more variables, there is a trend of narrowing down the resulting data's space due to the intersection of constraints. Therefore, this process should result in sets that are an equal or better assessment of the present state (than the sets generated by inputs). However, filtering might also result in zero samples (membership value of zero on all knowledge base's rows), which indicates not having "memory" of any similar occurrence. This zero-sample filter is an excellent indicator for potentially anomalous scenarios, primarily if we know that the knowledge base is statistically representative of the variables in the cell (its use is mentioned in 4.2.7).

This process also enables simple forecasting. Considering an offset (arbitrary number of rows/time forward) in the activations, the temporal similarity extraction and computation of outputs results in future values. So, cells may compute outputs related to the present or future. Nonetheless, this temporal shift might be difficult to achieve if data samples arrive at randomly spaced intervals of time, thus would only be straightforward for systems with time-consistent data acquisition.

Up next are practical examples of the temporal similarity extraction methodology. We consider that  $t = 0$  is associated with the present, and  $\mu$  represents membership functions.

**Self-similarity** The cell can perform intrinsic temporal similarity extraction with its knowledge base and input variables. Consider a cell characterized by two variables that are a function of time:  $x(t)$  and  $y(t)$ . Let us assume that, at a given instant, these are the cell's new inputs:

$$x(0) = 1.0 \rightarrow \mu_{x(0)}(x) = \begin{cases} 1, & x \in [0.9, 1.1] \\ 0, & x \notin [0.9, 1.1] \end{cases} \quad (4.3)$$

$$y(0) = 2.0 \rightarrow \mu_{y(0)}(y) = \begin{cases} 1, & y \in [1.8, 2.2] \\ 0, & y \notin [1.8, 2.2] \end{cases} \quad (4.4)$$

The membership functions  $\mu$  are generated considering that  $x$  and  $y$  are characterized by uniform and symmetrical uncertainty and that the received samples of  $x(0)$  and  $y(0)$  represent their median (classical set representation).

With the knowledge base represented in figure 4.5, the resulting activations will be 2023-01-01 00:02:00 and 2023-01-01 00:03:00. These are the timestamps of past instances where the cell's variables have values belonging to the set generated from new inputs. Now we can also infer that the actual values of  $x$  and  $y$  should reside in a set constrained by the filtered historical data ( $x(t)$  and  $y(t)$ ). Therefore, the outputs based on self-similarity extraction are:

$$x'(0) \in [0.90, 0.95] \quad (4.5)$$

$$y'(0) \in [1.80, 2.20] \quad (4.6)$$

timestamp	$x(t)$	$\mu_{x(0)}(t)$	$y(t)$	$\mu_{y(0)}(t)$	
2023-01-01 00:00:00	0.80	0	1.90	1	
2023-01-01 00:01:00	0.85	0	2.00	1	
2023-01-01 00:02:00	0.90	1	2.10	1	
2023-01-01 00:03:00	0.95	1	2.20	1	
2023-01-01 00:04:00	1.00	1	2.10	0	
2023-01-01 00:05:00	1.05	1	2.20	0	

Bounds of filtered knowledge

Figure 4.5: Visualization of a self-similarity extraction example.

These results confirm that we achieve outputs with less uncertainty only depending on this intrinsic process and private data.

**Mutual Similarity** Extending similarity extraction to multiple cells is a relatively trivial process. Suppose a cell has access to another’s activations at the same rolling timestamp and for a historical window that intersects its knowledge base. In that case, it can use that information to refine the intrinsic temporal similarity extraction result. Using sets, this can occur with a simple interval intersection of bounds. Although simple, this process has some tricky requirements, such as not allowing a time difference between the computation of membership values of different cells to avoid joining incoherent information and being able to intersect the activations of cells with different data acquisition periods.

When analyzing the resulting activations of the previous example, we noticed that using the timestamp of each activated row could be simultaneously more efficient and avoid synchronization issues. Firstly, representing activations by an interval of timestamps instead of the individual rows allows compressing consecutive samples into a single interval. This new representation may or may not improve efficiency since a knowledge base that filters into alternating activations doubles the memory requirement for this representation. However, an interval representation is suitable if inputs are uncertain and result in a large filter.

#### 4.2.5 Trust

After recognizing the potential for mutual time similarity extraction, we formulate a mechanism for understanding coherence between cells before mindlessly intersecting their activations. This unique characteristic of the CellTAN enables its usage without data privacy issues and not requiring data normalization. However, as seen before, the knowledge base must intersect for a time-based comparison to make sense, and it also yields the best results if outliers are absent (for anomaly detection). Besides cell comparison, activations also allow comparing variables by performing similarity extraction with a subset.

Time possesses inherent normalization properties, as all components equally experience it. This characteristic enables the network to compare systems without concerns about the specific variables involved. Consequently, this is why the proposed tool could be effectively applied beyond PV applications, showcasing its ability to generalize across various domains.

#### 4.2.5.1 Trust between Cells

A 2x2 contingency table represents the intersection of activations of two cells. In this table, the rows define the activation status of one cell (labeled as "Activated" and "Not Activated"), and the columns represent the activation status of the other cell (same labels).

		Cell 2	
		Activated	Not Activated
Cell 1	Activated	a	b
	Not Activated	c	d

Table 4.1: Contingency table representing the activation intersection of two cells.

Each element of the matrix represented in Table 4.1 ( $a, b, c$ , and  $d$ ) corresponds to a specific combination of activations, such as "Activated-Activated," "Activated-Not Activated," "Not Activated-Activated," and "Not Activated-Not Activated". They represent the frequency or count of observations falling into that specific combination. The main diagonal relates to the cells' agreeableness, while the secondary diagonal is the opposite. For this work, we decided to use the 'second' as the terms' unit, meaning that, for example,  $a$  is the total amount of seconds both cells are active.

The matrix representation of the activation intersection metrics resembles the statistical results of sampling a population with two binary categories. Therefore, a statistical test could be appropriate to calculate the association between these categories. The following section develops this topic.

#### 4.2.5.2 Statistical tests for measuring association

In statistics, numerous tests are available to measure the association between variables in contingency tables. One of the most known tests is Pearson's chi-square test (A.1.1), which assesses whether there is a significant association between the two variables. It compares the observed frequencies in the contingency table with the expected frequencies under an assumption of independence. If this test yields a statistically significant result, it suggests a non-random association between the activations of the two cells.

Another widely employed test is Fisher's exact test (A.1.2), as an alternative to the chi-squared, which is particularly useful for small sample sizes. It calculates the probability of obtaining the observed distribution of activations in the contingency table, also assuming independence between the variables. If this probability is sufficiently low, it implies that the association between the activations is unlikely to occur by chance. The odds ratio and relative risk can also quantify the

strength and direction of the association between two variables. The odds ratio would compare the odds of activation in one cell close to the other, while the relative risk compares the risk of activation in one cell to the risk in the other.

Considering that the desired outcome of activation comparison between cells is a normalized index that translates into how much they conform to each other (hence the term "trust"), there is a preference for inherently normalized tests, such as the  $\phi$  coefficient (A.1.4), contingency coefficient (A.1.5), and Theil's U (A.1.6). The goal is that this index represents how much the historical incidence of two cells' states match.

The  $\phi$  coefficient, also known as Matthews correlation coefficient, is a measure of association designed explicitly for 2x2 contingency tables, quite popular in machine learning for measuring the quality of binary classification. It ranges from -1 to 1, where 0 indicates no association, -1 represents a perfect negative association, and 1 represents a perfect positive association. The contingency coefficient extends this coefficient's usability, allowing its application in larger contingency tables, and ranges from 0 to 0.707 (no association to strong association). Finally, another relevant test is Theil's U . It accounts for the variables' mutual information and entropy based on information theory principles, providing a measure for the proportion of total entropy in one variable that the other can explain. It also ranges from 0 to 1.

#### 4.2.5.3 Proposed trust measurement method for cells

All the mentioned tests may offer different perspectives on the association between the activations of two cells. However, none are considered adequate for our application, so we propose a new method, defined by 4.7.

$$\chi_a^2 = \frac{(p_a - E(p_a))^2}{E(p_a)} \quad (4.7)$$

Where:

$$p_a = \frac{a}{a+b+c+d}$$

$$E(p_a) = \frac{(a+c) \times (a+b)}{(a+b+c+d)^2}$$

$p_a$  : Probability of both cells being active

$E(p_a)$  : Expected value for the probability of both cells being active

Based on Pearson's chi-squared test, this new approach focuses on the "a" component: the pure agreeableness between two cells. Its value ranges from 0 to 1 (no "trust" to complete "trust").

Some of the most relevant tests previously presented, namely the  $\phi$  coefficient, contingency coefficient, and Theil's U, are compared to this new method for comparison and validation in the following experiments.

#### 4.2.5.4 Experiments and validation

**Experiment n°1** The first trial simulates the variation of the  $d$  component and its effect on some association coefficients. It's expected that, as not activated time increases, the chance of the two cells intersecting activations lowers, and so their trust should increase. The starting point and used spaces of  $d$  are:

$$\begin{bmatrix} 30 & 10 \\ 10 & d \end{bmatrix}$$

$$d \in [0, 100], d \in [0, 3000]$$

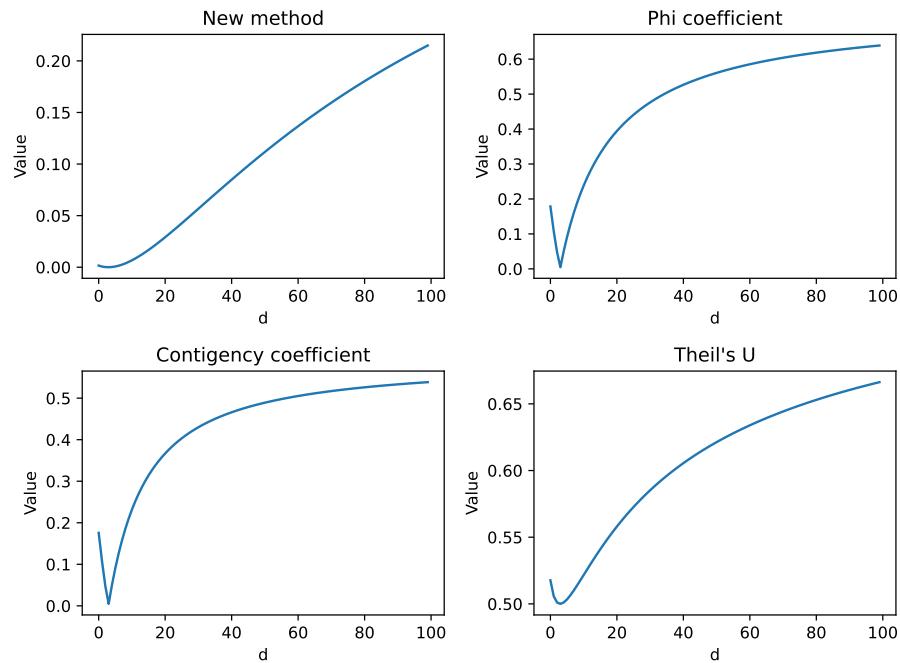


Figure 4.6: Trust measurement evolution for different methods with  $a = 30, b = 10, c = 10, d \in [0, 100]$ .

In figure 4.6, we can observe that all measurements increase with  $d$ , with a slight reduction at the start, although attenuated on the new method. Theil's U's minimum is 0.5, while the expectation for having a lower  $d$  (few deactivated intersecting periods) is a value closer to zero (as the other methods demonstrate). This scenario means either both cells have no history similar to current inputs or their inputs' values have significant uncertainty (remember the time extraction process from 4.2.4). Therefore, their "trust" should be minimal. Not having a reduction in the first values of  $d$ , showcased by the new method, is the intended behavior. Having equal activated and

deactivated intersections ( $a$  and  $d$ ) should not be significantly penalized versus having  $a \gg d$ , justified by the above reasoning.

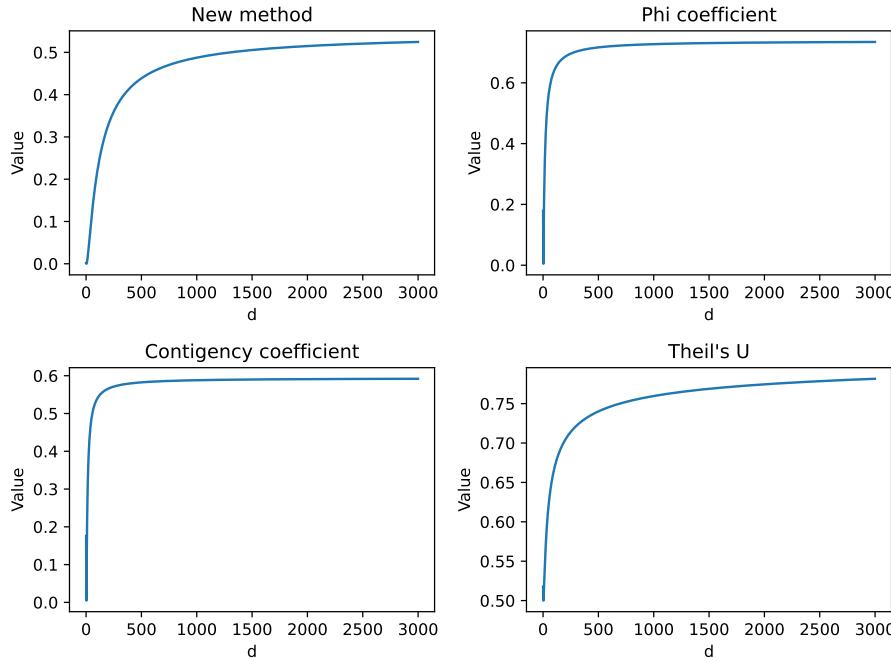


Figure 4.7: Trust measurement evolution for different methods with  $a = 30, b = 10, c = 10, d \in [0, 3000]$ .

Figure 4.7 extends the previous analysis, showcasing that all methods have asymptotical behavior, settling in a value greater than the start. All except Theil's U have asymptotes deemed acceptable for this test: since disagreement terms  $b$  and  $c$  are non-null and relatively close to  $a$ , we expect that the trust index does not reach a value near the max, remaining closer to a measure of "half trust".

**Experiment n°2** Now, we observe the impact of the disagreeing terms  $b$  and  $c$  in the same coefficients by varying one while maintaining all others fixed. The expectation is that an increase in disagreement should lead to zero trust.

$$\begin{bmatrix} 30 & b \\ 10 & 10 \end{bmatrix}$$

$$b \in [0, 1000]$$

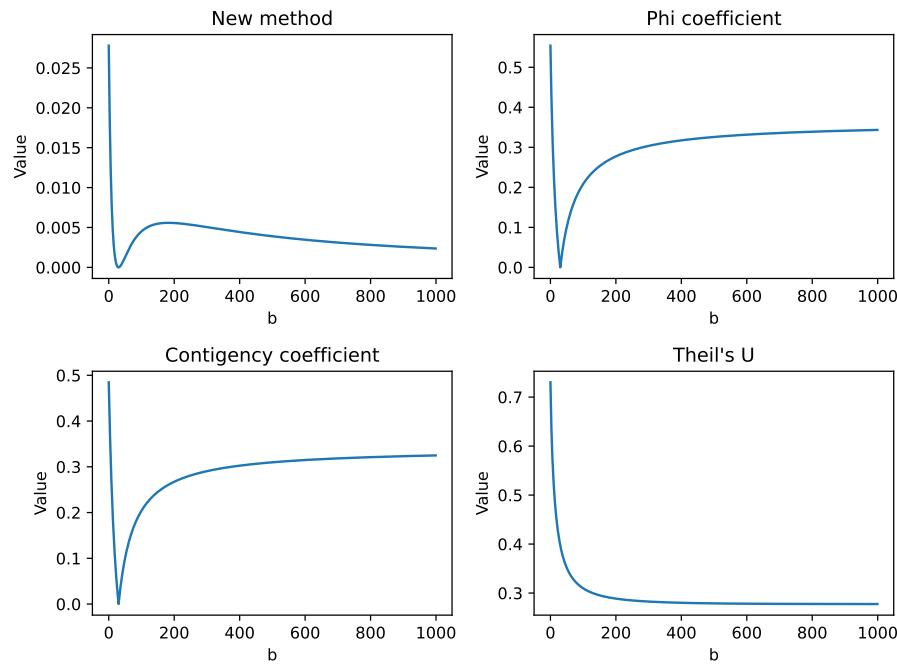


Figure 4.8: Trust measurement evolution for different methods with  $a = 30, c = 10, d = 10, b \in [0, 1000]$ .

In figure 4.8, it is clear that, for all methods, an increase in  $b$  leads to lower coefficients. All except the new algorithm tend towards similar asymptotes, near a value of 0.3. The new method meets the expected behavior.

$$\begin{bmatrix} 30 & 10 \\ c & 10 \end{bmatrix}$$

$$c \in [0, 1000]$$

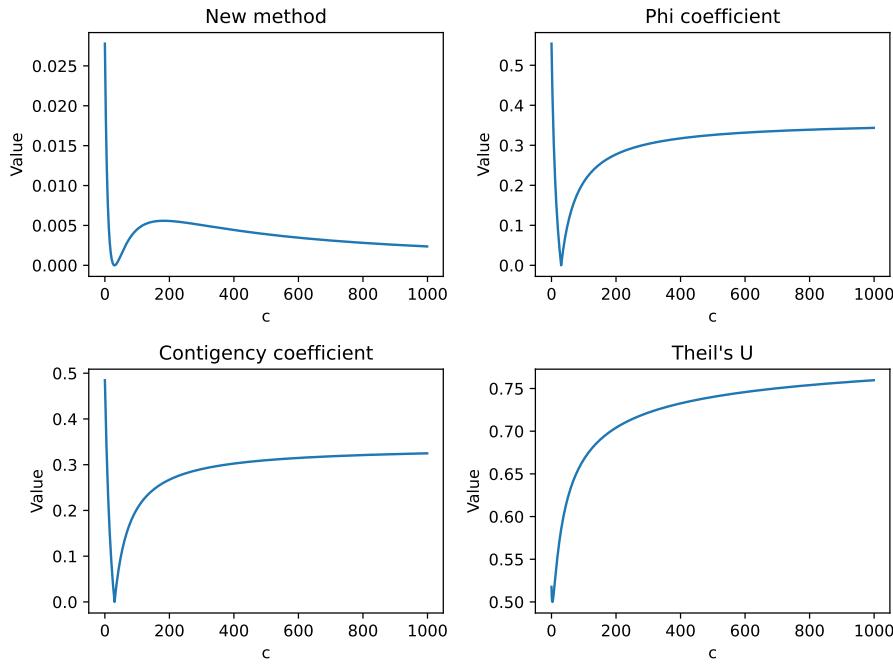


Figure 4.9: Trust measurement evolution for different methods with  $a = 30, b = 10, d = 10, c \in [0, 1000]$ .

With figure 4.9, we can observe the asymmetry of Theil's U coefficient, opposing to what figure 4.8 showed. It does not represent the trust between cell n°1 and cell n°2, but the trust of cell n°1 in cell n°2. The rest of the methods have symmetrical coefficients, thus presenting the same results.

**Experiment n°3** This last experiment showcases the effect of changing term  $a$ . When the cells share an increasing amount of activated time, assuming their trust should increase is trivial. However, reasoning that having a disproportionately ample activated time means there is more uncertainty on their current state, then trust should lower. The expectation is that, as  $a$  increases, there is also an increase in trust until a specific peak. It should decrease after reaching a global maximum and have an asymptote at zero.

$$\begin{bmatrix} a & 10 \\ 10 & 500 \end{bmatrix}$$

$$a \in [0, 1000]$$

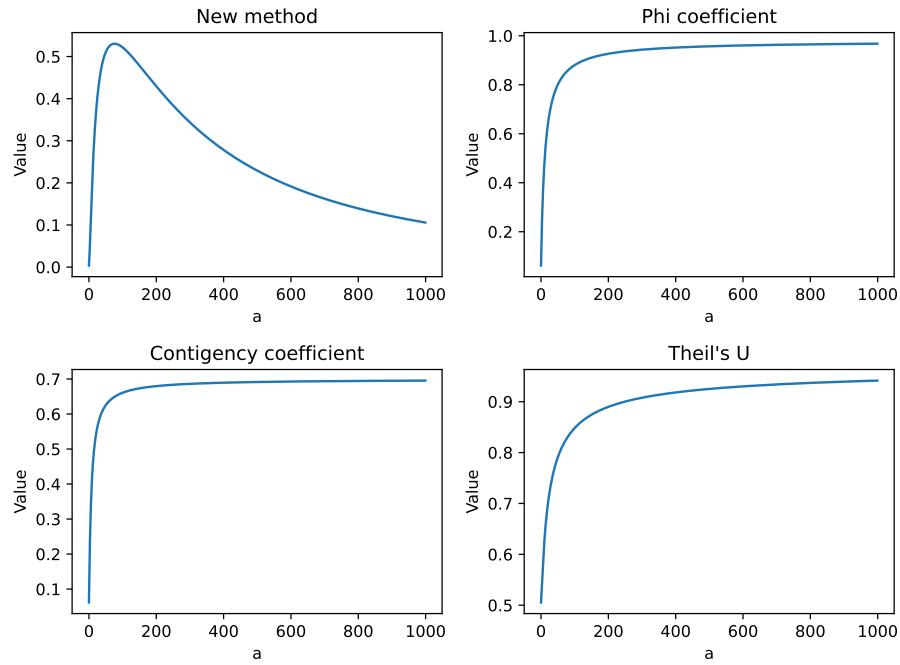


Figure 4.10: Trust measurement evolution for different methods with  $b = 30, c = 10, d = 10, a \in [0, 1000]$ .

Figure 4.10 shows that only the new method presents the previously described "goal" behavior. Otherwise, all other coefficients approach their max when  $a$  tends towards infinity.

With the three experiments presented, we confirm that the proposed approach for computing the trust index between cells is better suited, considering the requirements and assumptions.

#### 4.2.5.5 Trust between variables

Variables defined in the inputs/outputs of the cell can also be subject to activation comparison. Accordingly, we suggest simplifying this process compared to the inter-cell trust 4.8. Given that a cell may possess many variables, we restrict this computation to one per each, comparing the activations of a single variable with the intersection of all others.

$$Trust = 1 - \frac{TA_{others,exclusive}}{TA_{self,others} + TA_{self,exclusive} + TA_{others,exclusive}} \quad (4.8)$$

We utilize equation 4.8 to calculate the trust indicator for variables. The resulting value will be normalized and transparent, reflecting how much a variable was inactive while others were active.

#### 4.2.6 Connections

Mutual similarity extraction and trust calculation cannot occur in isolation. Therefore, connections are the essence of forming the CellTAN. These links can also have indicators for describing the cells' relationships, like the trust index, which is crucial for anomaly detection.

Each cell has a unique identifier, generated upon creation and independent of the name attributed by the owner. This ID serves to register cells and better manage the network. Since interconnections require an agent to keep track of these IDs and correctly redirect traffic, we introduce the **Hub** component. This central element provides global network visibility and makes connections easier to form and maintain. For a cell, a connection is no more than the identifier of the neighbor, which the hub can directly associate with a communication link.

On cell deployment, the system expects the owner and the CellTAN manager to manually link cells together by assigning each connection. However, the plan is to develop a mechanism that proposes new connections based on the neighbor's neighbors through the hub. This mechanism could function based on the strength of the relationship between cells and neighbors in common, recommending a direct link if both "trust" them similarly. As the last phrase hints, an indicator called "trust" is the criterion for evaluating these links. The following subsection defines the conception of this metric.

#### 4.2.7 Unconformities and State

We introduce a state system to expose the resulting unconformities found during cells' processes. These unconformities can categorized into intrinsic and extrinsic, based on the historical check of trust in the variables and neighbors, respectively, and during output computation.

Cells rely on aggregated trust values to assess the integrity and coherence of information. During specific time windows, these systems employ a verification process to identify unconformities in the aggregated trust values. Unconformities are flagged when the aggregated trust values fall below a certain threshold, which the owner predetermines. To establish a framework for this assessment, the owner specifies time bins upon cell creation, which define intervals evaluated from the current rolling timestamp to the past. Consequently, when inconsistencies are detected, detailed information is recorded to identify the elements that exhibit incoherence and the corresponding time bin in which the discrepancies occur. This approach allows for the characterization and analysis of "Intrinsic/Extrinsic time bin unconformities," thereby enabling the system to escalate the severity of unconformities from none to partial or total if some or all elements within a time bin exhibit incoherence.

In addition to the assessment methodology mentioned above, intrinsic filtering and the intersection of neighbor's activations play a crucial role in determining the overall unconformity state of the cell. Specifically, suppose intrinsic filtering yields a null domain. In that case, it indicates that the cell has encountered total intrinsic unconformity, suggesting a fundamental breakdown or lack of coherence with its knowledge base. Similarly, suppose the intersection of the neighbor's activations yields a null domain. In that case, it signifies total extrinsic unconformity, implying that the cell lacks trust in its external environment. These higher states of unconformity, whether intrinsic or extrinsic, serve as valuable indicators called the "Unconformity severity info." Considering this additional information, the system can escalate the severity level set by the time bin unconformities, providing a more comprehensive understanding of the cell's overall integrity and the extent of the identified unconformities.

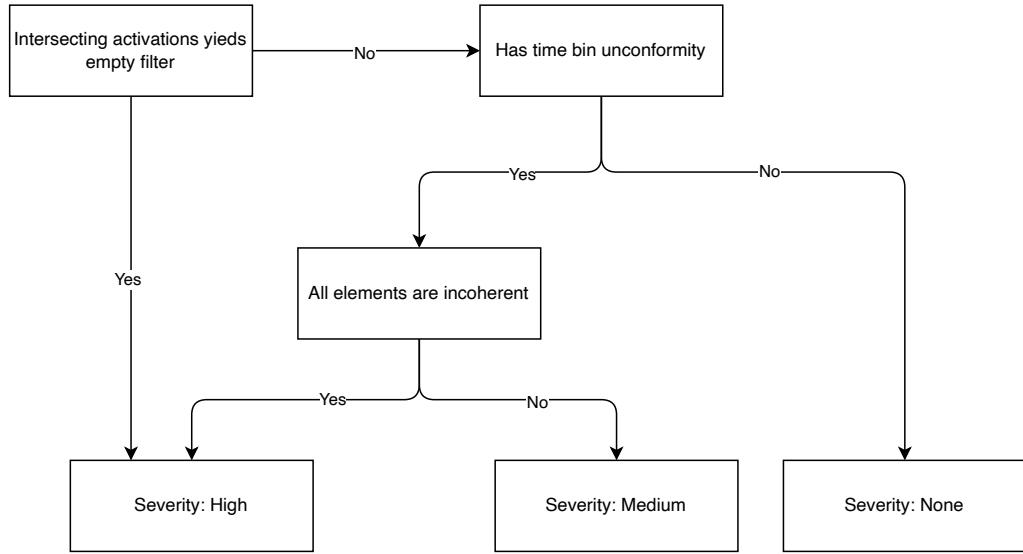


Figure 4.11: Algorithm for deciding unconformity severity (intrinsic and extrinsic).

Figure 4.11 sums up the algorithm of choosing unconformity severity. The severity level is associated with each type of unconformity since incoherence between neighbors and variables is not necessarily related. The first transition happens during output computation and the time bin unconformities is found during a historical trust assessment (processs seen in figure 4.2).

Though informative, the meaning of the unconformities and their severities is relatively broad, causing the cell to overlook specific situations related to its physical behavior. The following section addresses this, suggesting an application-oriented approach to extend observability.

#### 4.2.8 Application Plugins

The broad capabilities of the CellTAN make it less insightful when extracting specific information from the cells. Although its generalized design extends usability to different types of systems, it still requires additional processes to assess particular situations beyond what the state provides. During the literature review process, the specificity of algorithms is noticeable, and most techniques do not attempt to diverge from the PV application when it comes to fault detection and classification. This standard approach yields excellent results and could boost CellTAN's capabilities.

To solve the presented issue and dwell on the PV application theme of this work, we introduce the concept of **plugins**. Plugins extend the cell's core functionality, introduced to leverage specific system logic to identify patterns in its variables, just like classical algorithms. When binding them to a cell, they have a setup that checks if the required variables are present in its inputs. Then, after the intrinsic processes, they can run an algorithm that assesses what current situation the cell might be experiencing. This way, the CellTAN can still work with entities of different natures while these plugins work alongside, accessing private data and reporting sensible insights to the

owner. Nonetheless, they are confined to the cell's scope since permitting access to other cells' variables violates the privacy principle.

To illustrate this feature, let us imagine setting up a CellTAN representative of a solar farm with PV inverter cells. Although the tool does not "know" anything about inverters nor "cares" that the cells are of this type, the agent responsible for adding them to the network certainly does. Therefore, binding a PV-specific algorithm to assess particular inverter situations, such as malfunction, over(/under) voltage(/current), and performance degradation, is wise. This additional process makes the tool warn the owner of specific inverter faults besides possible mismatches in neighboring inverters and variables through the generalized cell procedures. The formulation of a PV plugin, used in a real case study of the CellTAN, happens in chapter 5.

### **4.3 The Hub**

The CellTAN tool is supposed to be easily accessible for different assets and owners in any given location. However, for privacy and security reasons, monitoring equipment and other smart devices (IoT, servers, etc.) in PV plants and the owner's database are usually protected from unwanted outside connections. Because of this, the CellTAN network owner can act as a proxy and be responsible for arranging the necessary connections between the equipment of different asset owners and the network. This way, all traffic is routed through its system, which solves the issues mentioned before but possibly introduces a bottleneck and affects availability. These downsides are inevitable for aggregating distributed systems and sharing information between otherwise reserved agents. Besides being a proxy, other responsibilities associated with the hub are:

- Provide network visibility: cells and connections;
- Receive events from cells;
- Cell connection proposal.

Adding this component permits visualizing all the cells registered in the system and their public data...

### **4.4 Expected Network Behavior**

In this section we can also notice how much the time decay mechanism of the inputs influences this process. If a cell is left unattended for a large enough amount of time, filtering with highly uncertain inputs leads to activations that cover most (or all) of its knowledge base. When this happens, outputs are equally uncertain and

## 4.5 Implementation

### 4.5.1 Code and Infrastructure

Materializing both the cell and the hub happened by coding Python modules. It was developed using a mix of the OOP (Object Oriented Programming) and FP (functional programming) paradigms and features a structure familiar with the descriptions in previous sections. Using Python for the implementation of the CellTAN comes with the following advantages:

- Easy to read and write code, requiring less syntax for complex operations compared to low-level languages;
- Extensive availability of libraries and tools;
- Deployment ease: does not have to compile, only needing an interpreter and dependencies to run;
- Big community support, with many resources publicly available online (e.g., documentation, tutorials, etc.).

Docker containers [45] are the infrastructure choice for deploying these modules (Cells and Hub). They allow running software as containerized applications, with all the necessary dependencies installed in an isolated environment. It acts as a separate system built for running the application instead of relying on the host's OS (Operating System) and running bare-metal. This execution strategy adds an isolation layer between the program and the host machine, increasing safety and making the "production" environment more predictable and stable. An overview of the technology stack utilized (and proposed) for the software products created in this work is pictured in appendix A.2.

### 4.5.2 Communication Protocol

In the context of ensuring proper data sharing between cells, the choice of communication protocol is a crucial aspect of establishing connections. For this work, we implemented a local communication mode for cells running on the same machine or within the same process and a remote communication model intended for cells distributed across different servers. The local communication mode is straightforward, assuming direct links between cells within the program. However, multiple protocol options regarding remote communication exist, including HTTP, HTTPS, WebSocket, and MQTT (excluding protocols for wireless proximity communication, as the CellTAN system assumes cells may be in different geographical locations). Initially, we implemented HTTP/HTTPS for remote communications due to its accessibility and ease of implementation. However, given its synchronous operation (based on request and response), we soon realized there might be more appropriate and robust choices for a distributed asynchronous system. Consequently, as future work, we should refactor the remote interface to utilize MQTT.

HTTP (Hypertext Transfer Protocol) and MQTT (Message Queuing Telemetry Transport) are communication protocols in different contexts. HTTP, a request-response protocol widely used in web applications, operates over TCP and follows a client-server model. It proves suitable when clients need to retrieve or send specific data to and from a server. HTTP is simple, widely supported, and compatible with browsers and standard web technologies. However, its stateless nature and lack of optimization for real-time communication can result in high overhead caused by frequent request-response cycles.

In contrast, MQTT is a lightweight publish-subscribe messaging protocol designed for efficient communication in distributed systems, especially within the Internet of Things (IoT). MQTT utilizes a broker-based architecture, where clients publish messages to topics and subscribe to receive messages from specific topics of interest. MQTT is highly scalable, bandwidth-efficient, and supports asynchronous messaging. It minimizes network and power consumption while providing reliable message delivery. However, implementing MQTT may require additional infrastructure, such as MQTT brokers. The publish-subscribe model of MQTT facilitates decoupled communication between system components, allowing for scalable and flexible systems. It is suitable for scenarios where devices or services exchange information in a distributed environment asynchronously. MQTT often outshines HTTP as the preferable choice for a distributed asynchronous system emphasizing real-time data exchange and efficient communication. Nevertheless, HTTP remains the more appropriate choice if the system primarily revolves around traditional web applications and request-response interactions. Ultimately, the scope of CellTAN requires the leverage of both protocols in different aspects of the tool: HTTP is more suitable for human-system interaction, and MQTT for system-system exchanges.

#### **4.5.3 Cell configuration and deployment**

Configuring cells can be done through configuration files (one per cell). They should contain the cell variables' definitions, database credentials, hub credentials, and all other parameters. Because of its simplicity, we chose the YAML serialization specification [46] to parse these configurations. A walkthrough of the cell configuration and expected file structure is present in appendix A.3

# Chapter 5

## CellTAN Application

### 5.1 Case study

The experiments validating CellTAN’s behavior incorporate two neighboring grid-tied string inverters from the same PV farm with common satellite data. Their only known characteristics are:

- Inverter one: 12.5kW nominal power, 14.4kW peak power. Installed January 1st, 2013.
- Inverter two: 15kW nominal power, 15.84kW peak power. Installed January 1st, 2013.

Variable	Source	Unit	Label
AC side power	Inverter (1 & 2)	W	ac_power
AC side current	Inverter (1 & 2)	A	ac_current
AC side voltage	Inverter (1 & 2)	V	ac_voltage
DC side power	Inverter (1 & 2)	W	dc_power
DC side current	Inverter (1 & 2)	A	dc_current
DC side voltage	Inverter (1 & 2)	V	dc_voltage
Global tilted irradiance	Satellite	W/m <sup>2</sup>	global_tilted_irradiance
Global horizontal irradiance	Satellite	W/m <sup>2</sup>	global_horizontal_irradiance
Cloud coverage	Satellite	%	cloud_coverage
Air temperature	Satellite	°C	temperature

Table 5.1: Available variables from two inverters and a satellite.

Table 5.1 represents the available variables and corresponding labels used to identify them in the cell’s inputs and graphs. These variables are sampled every 10 minutes from May 31, 2020, at 5:00 am to April 30, 2023, at 7:30 pm (although having some gaps). We utilized data from 2020 until the end of 2022 for the cells’ knowledge base, and any information from 2023 onwards is considered new and used for testing. Since there is no production at night, the database does not store values for this period. Not accounting for the night as missing samples, we have around 98% of data availability.

Analyzing and cleaning raw inverter and satellite data is essential to take full benefit of CellTAN's capabilities. As seen in its development, having a clean knowledge base contributes to correctly identifying anomalous situations. Therefore, the following sections focus on these two steps, contributing to understanding the anomalies' domain and frequency of occurrence.

### 5.1.1 Data analysis

Before data visualization, and regarding the variables in table 5.1, we eliminate those that will not benefit the CellTAN. We determined that AC side voltage is insignificant since the grid mandates it in a grid-tied inverter. We have decided to only use the measure of power instead of using the AC side current measure since, in conjunction with voltage, it provides the same information. To simplify things further, we do not need to consider the power on the DC if considering both the DC side current and voltage measures.

We examine all variables related to satellite data to determine which ones could be useful, not making any premature assumptions.

#### 5.1.1.1 Power

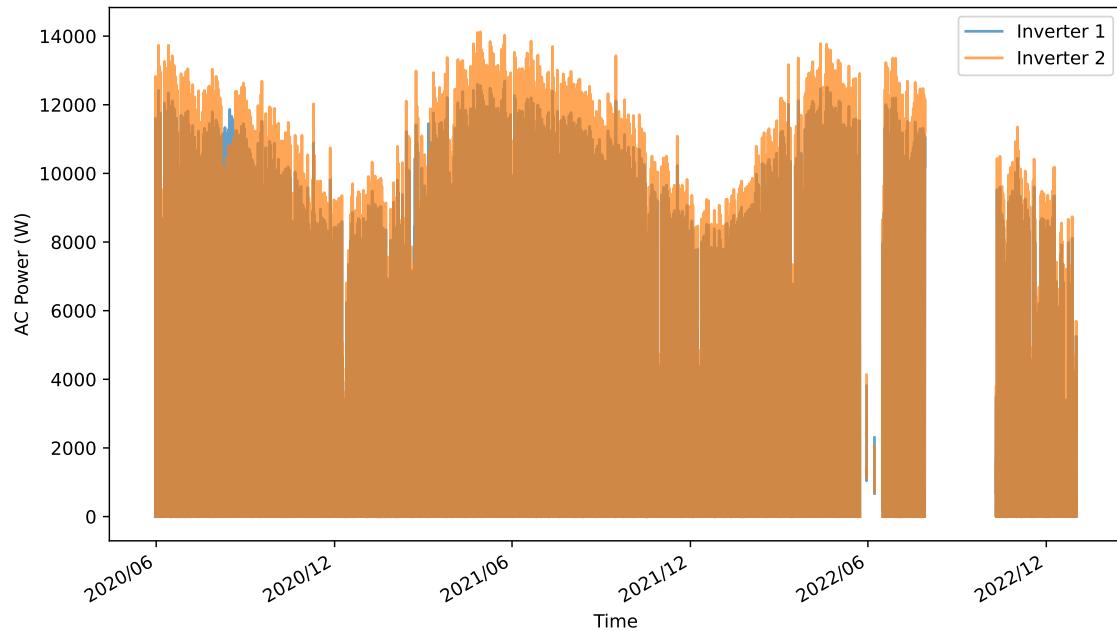


Figure 5.1: Inverter AC side power from 2020 to 2022, used for the knowledge base.

Figure 5.1 shows the power profile of the two studied inverters. Right away, we notice that the power of inverter two caps at around 14kW, while inverter one usually maxes at 12kW. This information is coherent with their ratings. We can notice two relatively large chunks of missing data, with the gaps occurring in mid to late 2022.

Figure 5.2 represents the power profile on the portion of data used for testing. When performing a closer inspection (with more zoom), we could hand-pick some fault occurrences in both

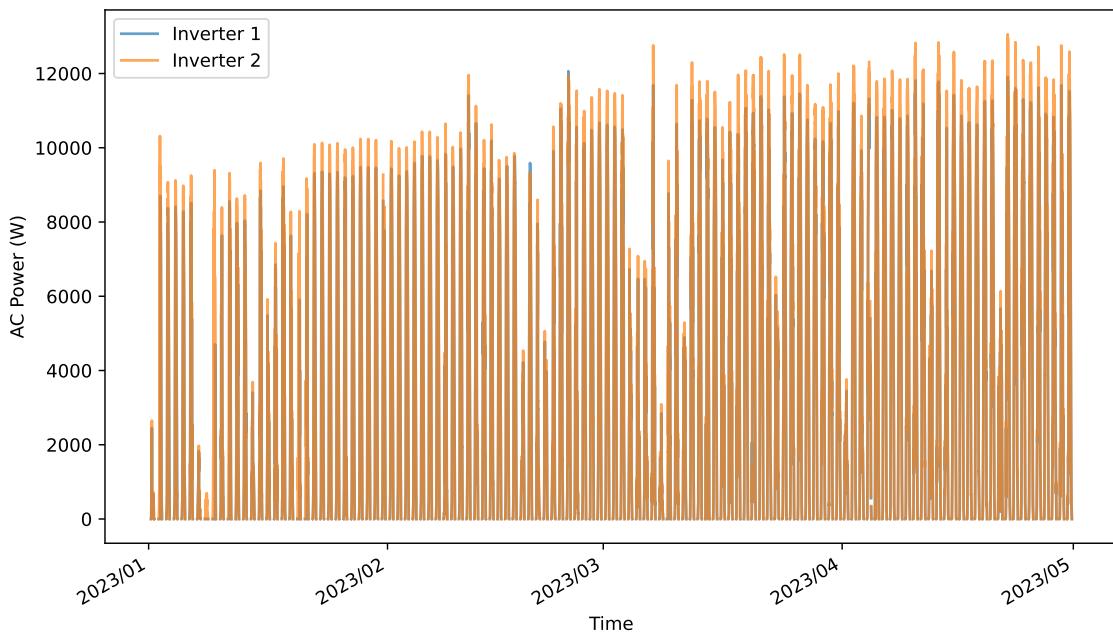


Figure 5.2: Inverter AC side power from 2023-01-01 to 2023-01-05, used for testing.

datasets, with the majority being one inverter off while the other continues regular operation. However, these will be more noticeable during different types of data analysis, such as pair plotting. Regardless, the cases that will matter are in the test data since these scenarios will not exist in the knowledge after cleaning. In Section 5.1.5, you will find a selection of carefully chosen scenarios.

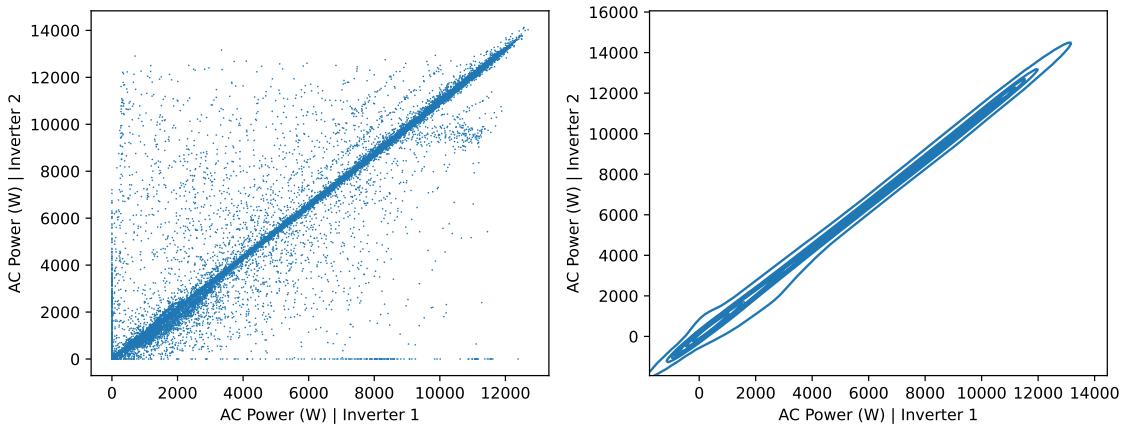


Figure 5.3: Pair plot of AC power from both inverters (2020 to 2022), using scatter (left) and KDE (Kernel Density Estimation) (right).

### 5.1.1.2 Voltage and Current

...

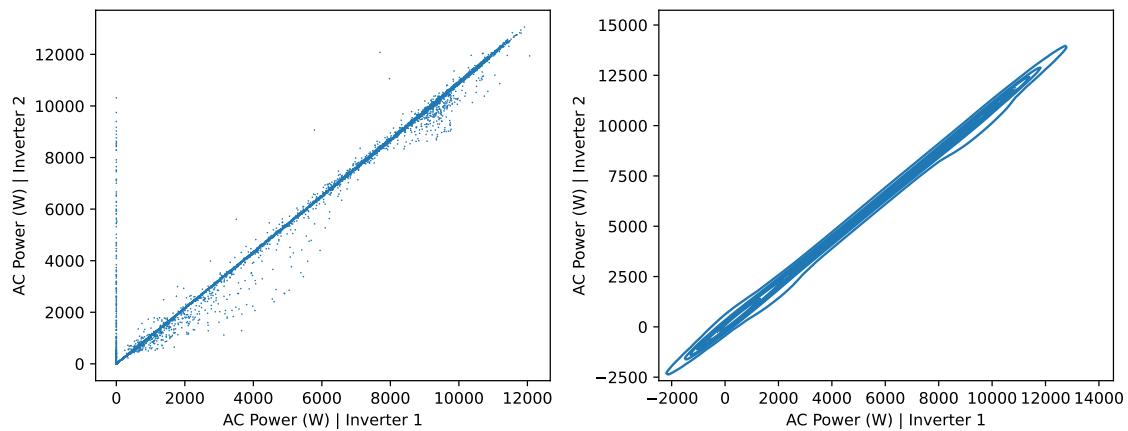


Figure 5.4: Pair plot of AC power from both inverters (2023), using scatter (left) and KDE (Kernel Density Estimation) (right).

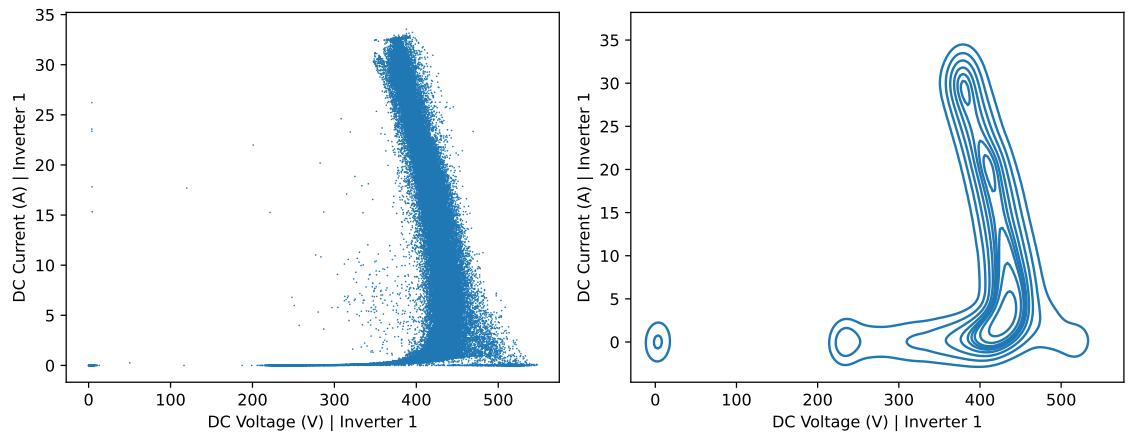


Figure 5.5: ...

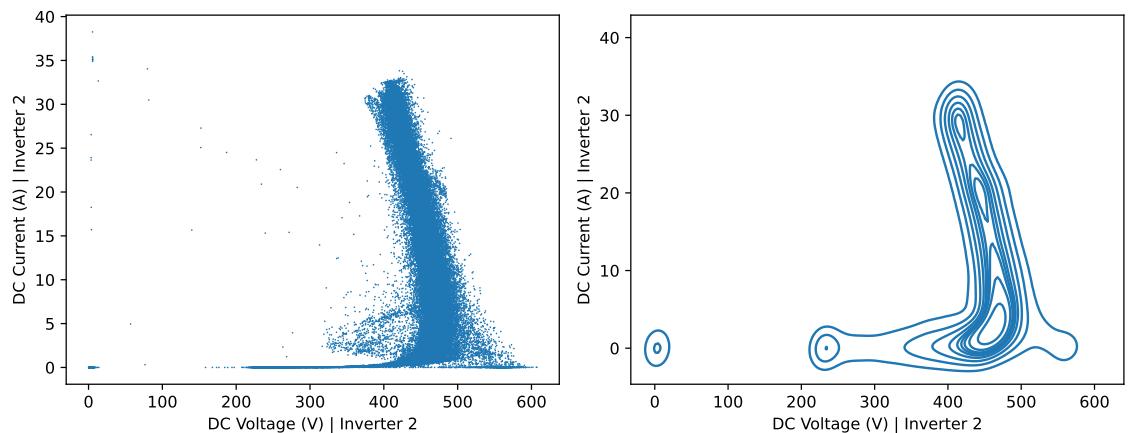


Figure 5.6: ...

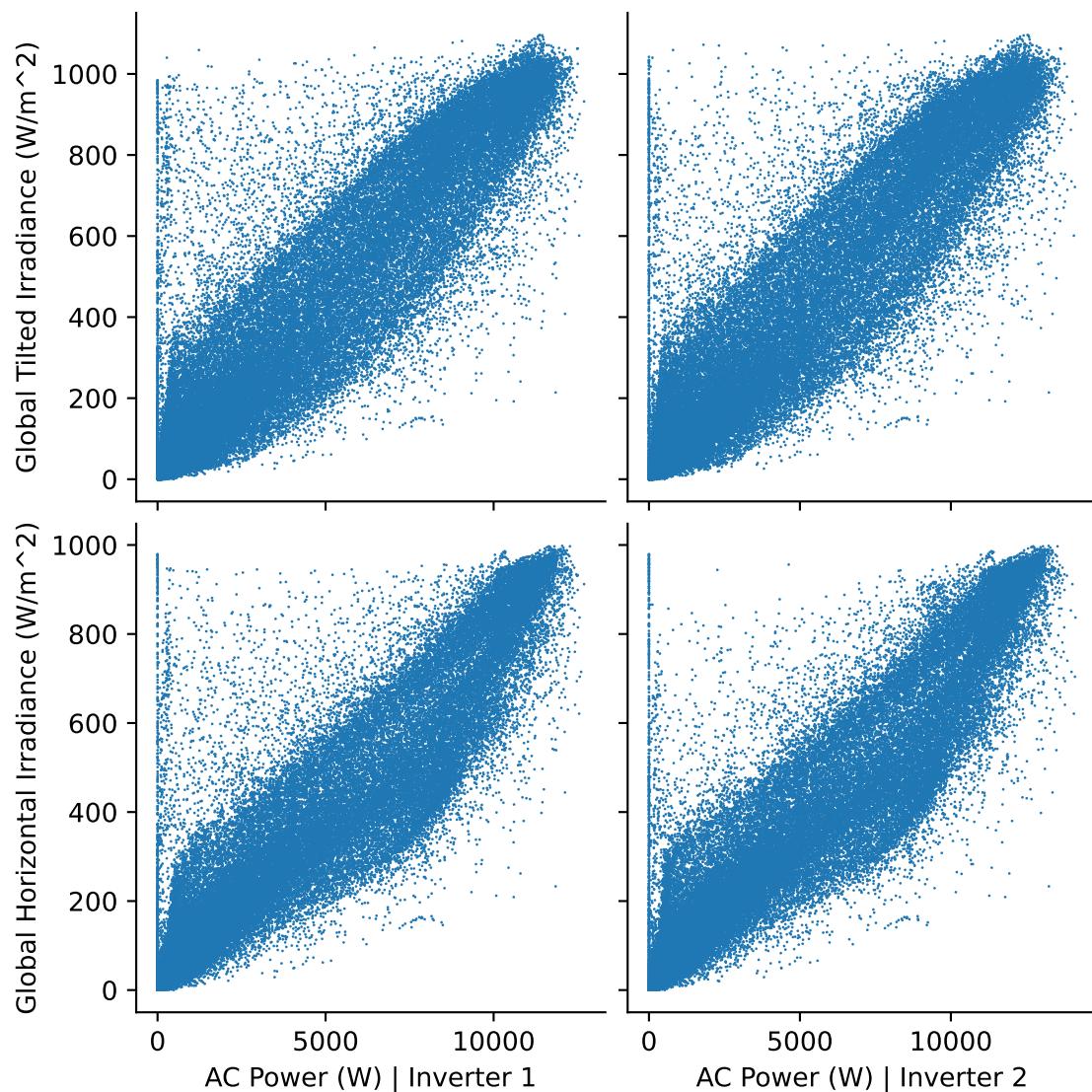


Figure 5.7: ...

### 5.1.1.3 Satellite

...

## 5.1.2 Data Cleaning

### 5.1.2.1 Power

...

### 5.1.2.2 Satellite

...

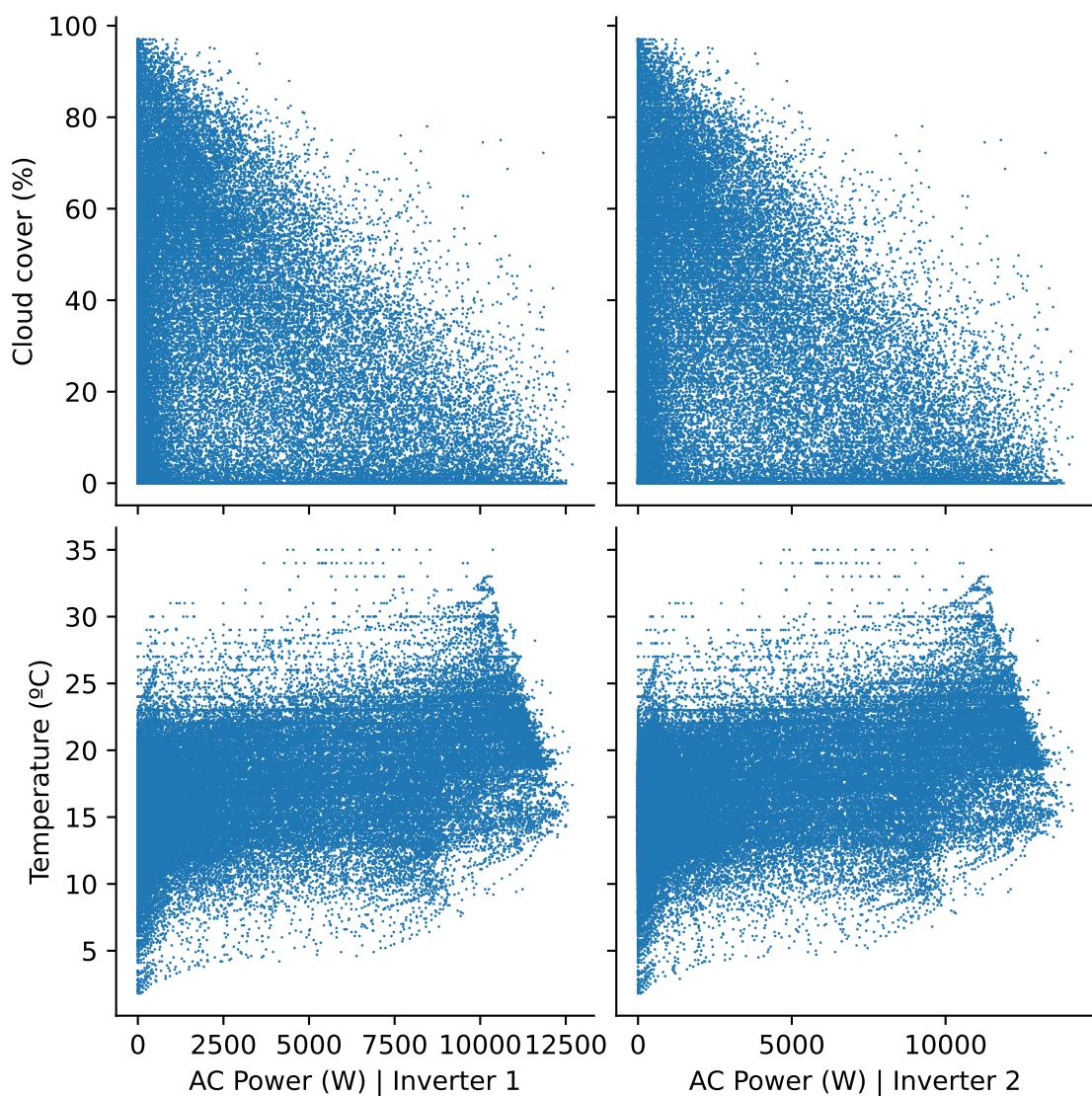


Figure 5.8: ...

### 5.1.2.3 Voltage and Current

...

### 5.1.3 Photovoltaic Plugin

...

### 5.1.4 CellTAN Configuration

...

**5.1.5 Simulation and Results**

...

**5.1.6 Scaling up**

...



# **Chapter 6**

## **Conclusion and Future Work**

...

### **6.1 Addressing the research questions**

...

### **6.2 Objectives reached**

...

### **6.3 Potential applications**

...

### **6.4 Future work**

...



# Appendix A

## CellTAN development

### A.1 Statistical tests for measure of association

#### A.1.1 Pearson's chi squared test

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \quad (\text{A.1})$$

where:

$\chi^2$  : Chi-squared statistic

$O_i$  : Observed frequency for category  $i$

$E_i$  : Expected frequency for category  $i$

$k$  : Number of categories or cells in the data

#### A.1.2 Fischer's exact test

$$p = \frac{\binom{a}{x} \binom{b}{y}}{\binom{N}{n}} \quad (\text{A.2})$$

where:

$p$  : p-value of the test

$a$  : Number of successes in group A

$b$  : Number of successes in group B

$x$  : Number of successes of interest in group A

$y$  : Number of successes of interest in group B

$N$  : Total number of observations

$n$  : Number of observations in group A

#### A.1.3 Odds ratio

$$OR = \frac{a \cdot d}{b \cdot c} \quad (\text{A.3})$$

where:

$OR$  : Odds ratio

$a$  : Number of successes in group A

$b$  : Number of failures in group A

$c$  : Number of successes in group B

$d$  : Number of failures in group B

#### A.1.4 Phi coefficient

$$\phi = \sqrt{\frac{\chi^2}{N}} \quad (\text{A.4})$$

where:

$\phi$  : Phi coefficient

$\chi^2$  : Chi-squared statistic

$N$  : Total number of observations

#### A.1.5 Contingency coefficient C

$$C = \sqrt{\frac{\chi^2}{N + \chi^2}} \quad (\text{A.5})$$

where:

$C$  : Contingency coefficient

$\chi^2$  : Chi-squared statistic

$N$  : Total number of observations

#### A.1.6 Theil's U

$$U(x|y) = \frac{H(x) - H(x|y)}{H(x)} \quad (\text{A.6})$$

Entropy of variable x:

$$H(x) = - \sum_{i=1}^n p(x_i) \log(p(x_i)) \quad (\text{A.7})$$

Conditional entropy of variable x given variable y:

$$H(x|y) = - \sum_{i=1}^n \sum_{j=1}^m p(x_i, y_j) \log \left( \frac{p(x_i, y_j)}{p(y_j)} \right) \quad (\text{A.8})$$

## A.2 Technology stack

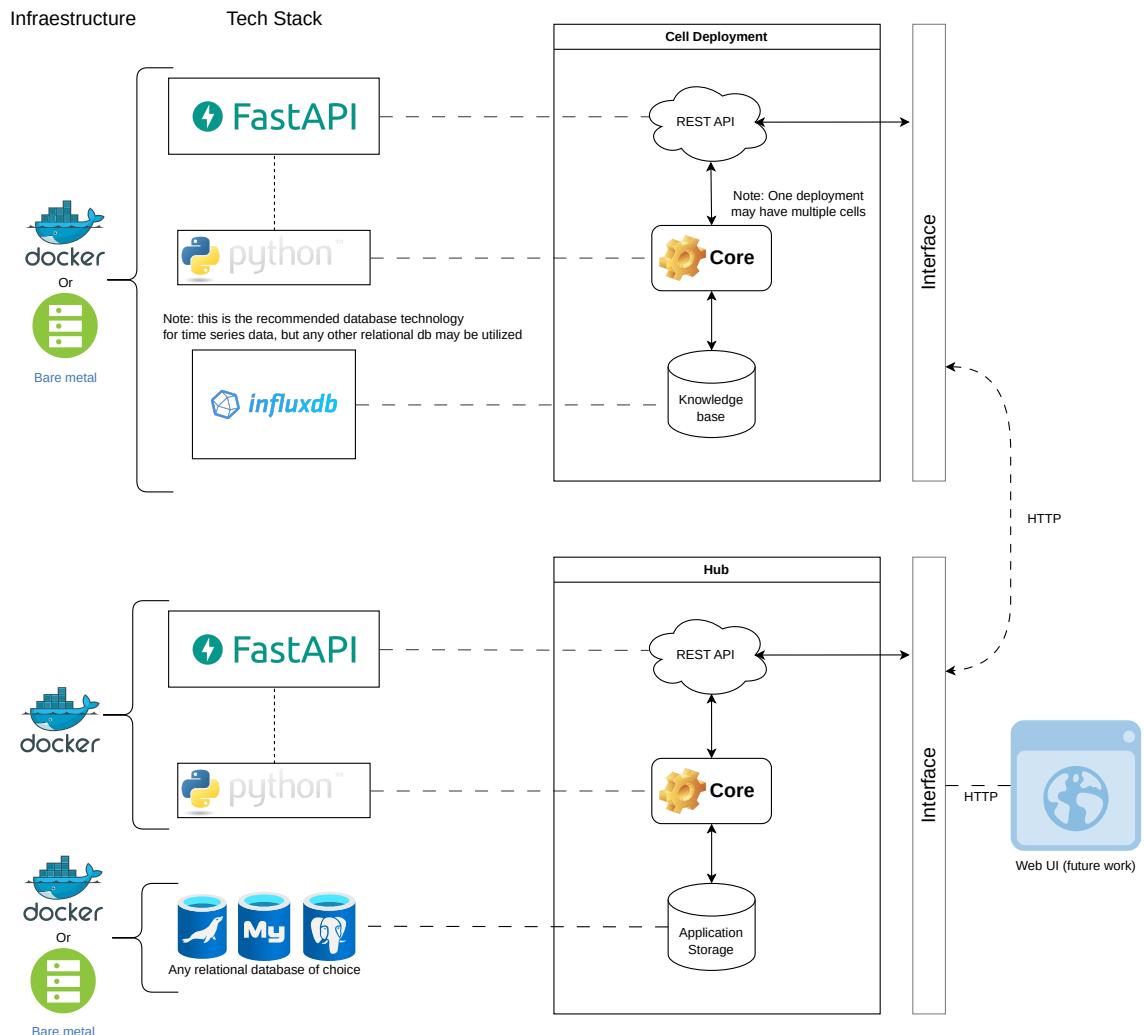


Figure A.1: Technology stack of the Cell and Hub of CellTAN.

## A.3 Cell configuration



# References

- [1] *Renewable Capacity Statistics 2021*. 2021.
- [2] A. Cabrera-Tobar, E. Bullich-Massagué, M. Aragüés-Peñalba, and O. Gomis-Bellmunt, “Topologies for large scale photovoltaic power plants,” 6 2016.
- [3] A. H. I. Mourad, H. Shareef, N. Ameen, A. H. Alhammadi, M. Iratni, and A. S. Alkaabi, “A state-of-the-art review: Solar trackers,” *2022 Advances in Science and Engineering Technology International Conferences, ASET 2022*, 2022.
- [4] A. Mellit and S. Kalogirou, “*Handbook of Artificial Intelligence Techniques in Photovoltaic Systems*”. Academic Press, 2022. <https://www.sciencedirect.com/book/9780128206416/handbook-of-artificial-intelligence-techniques-in-photovoltaic-systems>.
- [5] “Iec 61724-1:2021 rlv | iec webstore.” <https://webstore.iec.ch/publication/70170>. Accessed: 2023-01-07.
- [6] D. S. Pillai and N. Rajasekar, “A comprehensive review on protection challenges and fault diagnosis in pv systems,” 8 2018.
- [7] M. K. Alam, F. Khan, J. Johnson, and J. Flicker, “A comprehensive review of catastrophic faults in pv arrays: Types, detection, and mitigation techniques,” 5 2015.
- [8] H.-Y. Li, *Assessment of Photovoltaic Module Failures in the Field*. 05 2017.
- [9] Y. Y. Hong and R. A. Pula, “Methods of photovoltaic fault detection and classification: A review,” 11 2022.
- [10] A. Livera, M. Theristis, G. Makrides, and G. E. Georgiou, “Recent advances in failure diagnosis techniques based on performance data analysis for grid-connected photovoltaic systems,” 4 2019.
- [11] A. Stepanov, A. Sokolovs, and L. Dzelzkaleja, “Solar tracker supervisory system,” *2014 55th International Scientific Conference on Power and Electrical Engineering of Riga Technical University, RTUCON 2014*, pp. 79–83, 12 2014.

- [12] T. G. Amaral, V. F. Pires, and A. J. Pires, "Fault detection in pv tracking systems using an image processing algorithm based on pca," *Energies* 2021, Vol. 14, Page 7278, vol. 14, p. 7278, 11 2021.
- [13] H. Braun, S. T. Buddha, V. Krishnan, C. Tepedelenlioglu, A. Spanias, T. Takehara, T. Yeider, M. Banavar, and S. Takada, "Signal processing for solar array monitoring, fault detection, and optimization," 2011. summary:<br/>chapter 2 describes physics and electrical behaviour of the PV cell<br/>chapter 3 has description of the most common faults<br/>chapter 4 has statistical signal processing-based techniques for determining the presence<br/>of a fault.<br/><br/>chapter 2:<br/>"The most common metric comparing measured module performance to predicted behavior<br/>is the performance ratio"<br/>Modelling performance:<br/>Sandia Model: very accurate (1testing<br/>The five-parameter model: it utilizes only data provided by the manufacturer at standard test conditions and does not require additional measurements to derive parameters.
- [14] F. Grimaccia, S. Leva, A. Dolara, and M. Aghaei, "'survey on pv modules' common faults after an o&m flight extensive campaign over different plants in italy", *IEEE Journal of Photovoltaics*, vol. 7, pp. 810–816, 5 2017.
- [15] "Fire safety for pv systems - sunny. sma corporate blog." <https://www.sma-sunny.com/en/fire-safety-for-pv-systems/>. Accessed: 2023-01-07.
- [16] M. Kumar and D. V. S. K. Rao, "Modelling and parameter estimation of solar cell using genetic algorithm," *2019 International Conference on Intelligent Computing and Control Systems, ICCS 2019*, pp. 383–387, 5 2019.
- [17] R. Godina, E. M. Rodrigues, E. Pouresmaeil, and J. P. Catalão, "Simulation study of a photovoltaic cell with increasing levels of model complexity," *Conference Proceedings - 2017 17th IEEE International Conference on Environment and Electrical Engineering and 2017 1st IEEE Industrial and Commercial Power Systems Europe, EEEIC / I and CPS Europe 2017*, 7 2017.
- [18] N. R. Prasad, S. Almanza-Garcia, and T. T. Lu, "Anomaly detection," *ACM Computing Surveys (CSUR)*, vol. 14, pp. 1–22, 7 2009.
- [19] M. Ilas and C. Ilas, "Towards real-time and real-life image classification and detection using cnn: a review of practical applications requirements, algorithms, hardware and current trends," pp. 225–233, 10 2020.
- [20] S. Buddha, H. Braun, V. Krishnan, C. Tepedelenlioglu, A. Spanias, T. Yeider, and T. Takehara, "Signal processing for photovoltaic applications," *2012 IEEE International Conference on Emerging Signal Processing Applications, ESPA 2012 - Proceedings*, pp. 115–118, 2012.

- [21] Y. Zhao, F. Balboni, T. Arnaud, J. Mosesian, R. Ball, and B. Lehman, “Fault experiments in a commercial-scale pv laboratory and fault detection using local outlier factor,” *2014 IEEE 40th Photovoltaic Specialist Conference, PVSC 2014*, pp. 3398–3403, 10 2014.
- [22] S. Vergura, G. Acciani, V. Amoruso, and G. Patrono, “Inferential statistics for monitoring and fault forecasting of pv plants,” *IEEE International Symposium on Industrial Electronics*, pp. 2414–2419, 2008.
- [23] H. Iles and Y. Mahmoud, “Power based fault detection method for pv arrays,” *IECON Proceedings (Industrial Electronics Conference)*, vol. 2021-October, 10 2021.
- [24] J. Fan, S. Rao, G. Muniraju, C. Tepedelenlioglu, and A. Spanias, “Fault classification in photovoltaic arrays using graph signal processing,” *Proceedings - 2020 IEEE Conference on Industrial Cyberphysical Systems, ICPS 2020*, pp. 315–319, 6 2020. “We propose here a graph signal processing based semi-supervised learning technique, which achieves good performance in fault classification with relatively limited data.”  
Classified faults: standard test conditions (STC), shaded modules, degraded modules, soiled modules, and short circuit conditions.  
Also uses PVWatts dataset  
MOST SIMILAR WITH THE THESIS.
- [25] A. P. Dobos, “Pvwatts version 1 technical reference,” 2013.
- [26] S. Rao, D. Ramirez, H. Braun, J. Lee, C. Tepedelenlioglu, E. Kyriakides, D. Srinivasan, J. Frye, S. Koizumi, Y. Morimoto, and A. Spanias, “An 18 kw solar array research facility for fault detection experiments,” *Proceedings of the 18th Mediterranean Electrotechnical Conference: Intelligent and Efficient Technologies and Services for the Citizen, MELECON 2016*, 6 2016.
- [27] S. Katoch, G. Muniraju, S. Rao, A. Spanias, P. Turaga, C. Tepedelenlioglu, M. Banavar, and D. Srinivasan, “Shading prediction, fault detection, and consensus estimation for solar array control,” *Proceedings - 2018 IEEE Industrial Cyber-Physical Systems, ICPS 2018*, pp. 217–222, 6 2018. They have sensors per each panel, not realistic for utility scale pv plants  
good information about consensus in graphs  
consensus is used for determining average values along the entire array, based on each sensor’s measurements.
- [28] A. Kumari, A. Shekhar, and M. S. Kumar, “An artificial neural network-based fault detection technique for pv array,” *2022 2nd International Conference on Emerging Frontiers in Electrical and Electronic Technologies, ICEFEET 2022*, 2022. Uses ANN for classification  
Classified Faults: short circuit faults and hot spot faults  
Inputs:  $dI/dt$  and  $dV/dt$   
Results: 98.4 accuracy  
Pros: is that it can use the  $dI/dt$  and  $dV/dt$  as network inputs to assess faults, yields high accuracy.  
Cons: uses high frequency data from current and voltage of the panels, sampled every few milliseconds.
- [29] F. Aziz, A. U. Haq, S. Ahmad, Y. Mahmoud, M. Jalal, and U. Ali, “A novel convolutional neural network-based approach for fault classification in photovoltaic arrays,” *IEEE Access*,

- vol. 8, pp. 41889–41904, 2020. “utilizes deep two-dimensional (2-D) Convolutional Neural Networks (CNN) to extract features from 2-D scalograms generated from PV system data in order to effectively detect and classify PV system faults.”  
 “A survey study conducted in 2010 showed that such faults can reduce the generated power of photovoltaic systems annually by about 18.9to noice. Heavy deep learning model.
- [30] S. Rao, A. Spanias, and C. Tepedelenlioglu, “Solar array fault detection using neural networks,” *Proceedings - 2019 IEEE International Conference on Industrial Cyber Physical Systems, ICPS 2019*, pp. 196–200, 5 2019. Covers a lot of faults: 8 in total.
- [31] S. Rao, G. Muniraju, C. Tepedelenlioglu, D. Srinivasan, G. Tamizhmani, and A. Spanias, “Dropout and pruned neural networks for fault classification in photovoltaic arrays,” *IEEE Access*, vol. 9, pp. 120034–120042, 2021. Uses an autoencoder machine learning framework  
 Emphasizes the performance  
 Autoencoder detects the fault, and neural network classifies it  
 “We consider the approach of fault detection and classification by monitoring the electrical signals such as maximum power point tracking (MPPT) parameters”  
 They use the THE PVWatts DATASET  
 BEST ARTICLE SO FAR!
- [32] H. Kilic, B. Khaki, B. Gumus, M. Yilmaz, and P. Palensky, “Fault detection in photovoltaic arrays via sparse representation classifier,” *IEEE International Symposium on Industrial Electronics*, vol. 2020-June, pp. 1015–1021, 6 2020. Detect faults: DC short circuit faults of PV array. line to line and line to ground  
 Has a lot of sources and compares a lot of them!  
 “The aim of the proposed method is to detect the faults under low-mismatch, low-irradiance and high-impedance conditions.”  
 “The common faults in PV systems are short and open circuit faults as well as panel mismatch and module failures”  
 Relatively light algorithm that doesn’t require tuning.
- [33] G. Uehara, S. Rao, M. Dobson, C. Tepedelenlioglu, and A. Spanias, “Quantum neural network parameter estimation for photovoltaic fault detection,” *IISA 2021 - 12th International Conference on Information, Intelligence, Systems and Applications*, 7 2021. QNN are still in early stages of development  
 Not as good as classical Neural Network  
 Promising for the future.
- [34] I. Høiaas, K. Grujic, A. G. Imenes, I. Burud, E. Olsen, and N. Belbachir, “Inspection and condition monitoring of large-scale photovoltaic power plants: A review of imaging technologies,” 6 2022.
- [35] A. K. V. de Oliveira, M. Aghaei, and R. Rüther, “Automatic inspection of photovoltaic power plants using aerial infrared thermography: A review,” 3 2022.
- [36] X. Li, Q. Yang, Z. Lou, and W. Yan, “Deep learning based module defect analysis for large-scale photovoltaic farms,” *IEEE Transactions on Energy Conversion*, vol. PP, pp. 1–1, 10 2018.

- [37] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Communications of the ACM*, vol. 60, pp. 84–90, 6 2012.
- [38] M. Hajij, K. Istvan, and G. Zamzmi, “Cell complex neural networks,” 2020.
- [39] “About us - enlitia.” <https://www.enlitia.com/about>. Accessed: 2023-06-03.
- [40] “Python.” <https://www.python.org/>. Accessed: 2023-01-09.
- [41] “Tensorflow.” <https://www.tensorflow.org/>. Accessed: 2023-01-09.
- [42] “Scikit-learn.” <https://scikit-learn.org/stable/>. Accessed: 2023-01-09.
- [43] “Scipy.” <https://scipy.org/>. Accessed: 2023-01-09.
- [44] D. li Zhang, C. Guo, and D. Chen, “On generalized fuzzy numbers,” *Iranian Journal of Fuzzy Systems*, vol. 16, 2019.
- [45] “Docker.” <https://www.docker.com/>. Accessed: 2023-05-15.
- [46] “Yaml.” <https://yaml.org/>. Accessed: 2023-05-15.