

UWLID: 21500835

Applications of AI

Dr Nasser Matoorian

Abstract

In today's digital landscape, physical servers form the backbone of data centres, IT infrastructures and small businesses, playing a crucial role in ensuring seamless data management, security, and processing (Kumar, U. 2018). As reliance on these servers grows, maintaining their optimal performance becomes vital for business continuity. Traditional maintenance approaches, including reactive and preventive strategies, struggle to anticipate potential failures, leading to unexpected downtime and costly repairs. This project explores the transformative potential of predictive maintenance, harnessing the power of artificial intelligence (AI) and machine learning (ML) for server upkeep. By leveraging TensorFlow, a cutting-edge ML framework, the project aims to design and develop a predictive maintenance model that analyses server sensor data, such as CPU temperature and fan speed, to predict maintenance needs accurately. This predictive approach aims to reduce downtime, extend server lifespan, and optimise operational efficiency. The final model, implemented within a Jupyter notebook, will provide real-time insights into server health, guiding proactive maintenance decisions and revolutionising how server management is approached in enterprise IT infrastructures.

The sensor data in question is not historical, this project will use data collected manually from physical servers discussing how it is collected, organised, manipulated whilst providing justifications as i proceed.

In addition, there will be accompanying youtube videos, made by myself, to further add context to the following chapters: "Servers", "Data Collection" and "Data Organisation"

Abstract.....	1
Honourable Mentions.....	3
Background.....	4
Introduction.....	6
Aims and Objectives.....	8
Aim:.....	8
Objectives:.....	8
Data Collection and Preprocessing:.....	8
Exploratory Data Analysis:.....	8
Model Design and Development:.....	8
Model Training and Validation:.....	9
Integration and Deployment:.....	9
Analysis of Practical Impact:.....	9
Introduction into Essential ML Tools.....	10
TensorFlow.....	10
NumPy.....	10
Pandas.....	11
Matplotlib.....	11
Scikit-learn.....	11
Flask.....	12
Docker.....	12
Google Cloud AI Platform.....	12
GitHub.....	12
Servers.....	14
Central Processing Unit (CPU).....	16
Fans.....	17
RAID Controller (Redundant Array of Independent Disks).....	18
RAM (Random Access Memory).....	18
Network Interface Controller.....	19
Hard Disk Drive.....	19
PSU (Power Supply Unit).....	20
Overall.....	21
Data Collection.....	23
Data Organisation.....	27
Supervised Learning.....	31
Python Script.....	33
Initial_Model.....	34

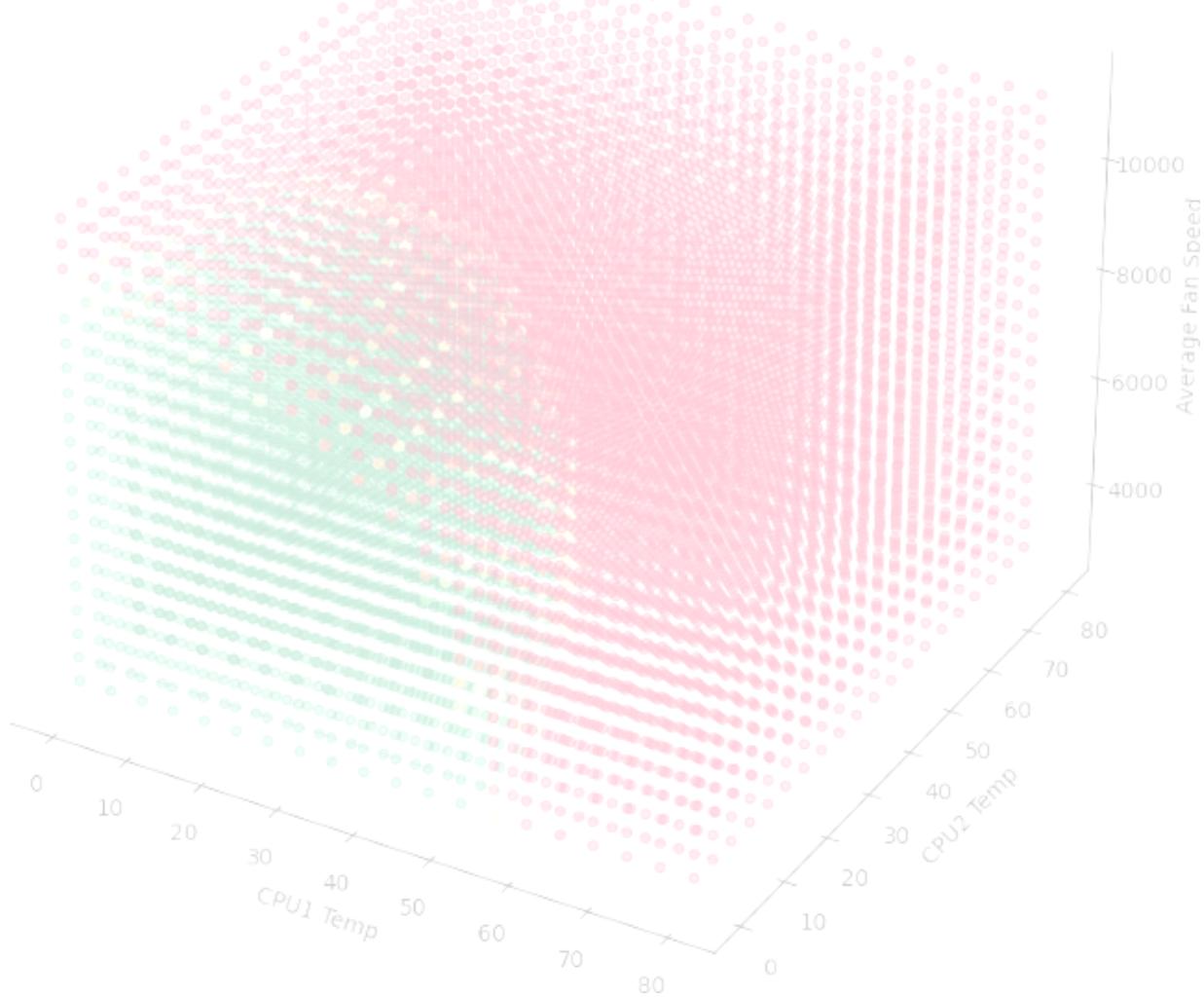
Honourable Mentions

Nathan Hollowbread - Warehouse Manager (IT Business Group Ltd)

Greg Seridarian - Managing Director (IT Business Group Ltd)

Dr Nasser Matoorian - Head of Computer Science (University of West London)

Dr Massoud Zolghani - Lecturer (University of West London)



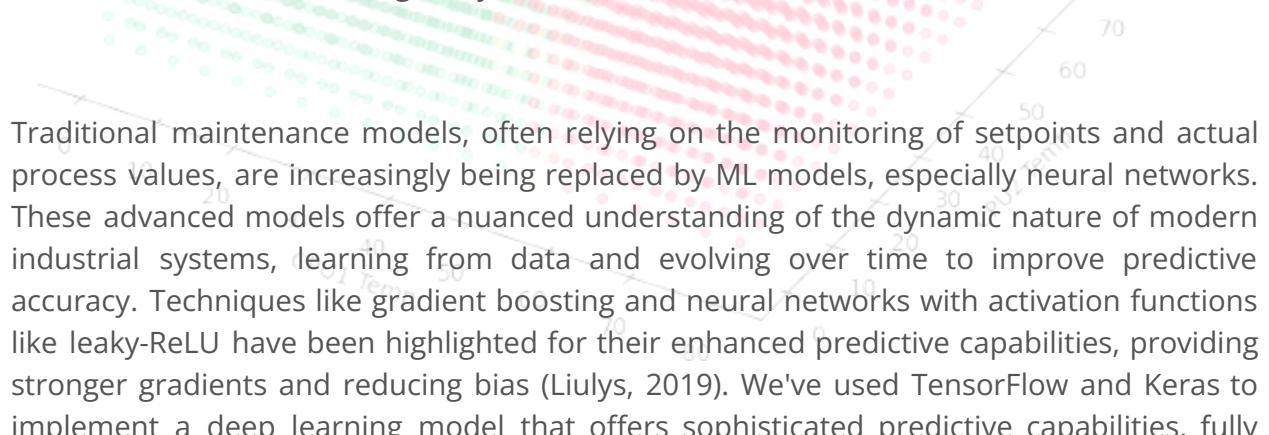
Background

The advancement of predictive maintenance (PdM) strategies, as part of the transformative wave of Industry 4.0, is deeply intertwined with the rapid development in machine learning (ML) and artificial intelligence (AI). This evolution marks a significant departure from traditional maintenance practices, moving towards more sophisticated, data-driven approaches.



The integration of machine learning and web frameworks into PdM has enabled predictive models to be packaged and deployed with ease, enhancing real-time monitoring capabilities. By utilising tools like TensorFlow, Flask, Docker, and Google Cloud AI Platform, the predictive maintenance landscape is transformed to foster real-time insights and enhanced decision-making capabilities.

The integration of the Internet of Things (IoT) into industrial systems has been a game-changer, facilitating the unification of various devices into a cohesive system. This integration allows for predictive and preventative maintenance strategies to use advanced ML algorithms, fundamentally altering the maintenance landscape in industrial settings (Liulys, 2019). In our project, IoT's role in gathering and analysing vast amounts of data is pivotal in enhancing equipment efficiency and reliability, aligning with Industry 4.0's ethos of interconnected and intelligent systems.



Traditional maintenance models, often relying on the monitoring of setpoints and actual process values, are increasingly being replaced by ML models, especially neural networks. These advanced models offer a nuanced understanding of the dynamic nature of modern industrial systems, learning from data and evolving over time to improve predictive accuracy. Techniques like gradient boosting and neural networks with activation functions like leaky-ReLU have been highlighted for their enhanced predictive capabilities, providing stronger gradients and reducing bias (Liulys, 2019). We've used TensorFlow and Keras to implement a deep learning model that offers sophisticated predictive capabilities, fully aligned with the evolution of PdM strategies.

The transition from simple Run-to-Failure (R2F) methods to more complex and efficient PdM systems is well documented. PdM systems now predict pending failures using

historical data, defined health factors, and statistical inference methods. This shift, driven by increasing data availability and the capabilities of modern hardware and algorithms, underscores the growing effectiveness of ML solutions in maintenance management (Susto et al., 2015). By containerizing the model using Docker, we ensured that our predictive model is portable and easily deployable across different environments.

The research trend in ML has shifted to more complex models, such as ensemble methods and deep learning, due to their higher accuracy in managing large datasets. The rise of deep learning, in particular, owes much to advancements in computing power, notably the evolution of GPUs. These developments have made deep learning a prominent research focus, with its ability to handle the complexities of large-scale industrial data (Research Paper, 2020). In our project, the TensorFlow library was crucial for building and training neural networks capable of deciphering patterns in server data, enabling early detection of potential issues.

Industry 4.0, characterised by cyber-physical systems and the industrial internet of things, integrates software, sensors, and intelligent control units. This integration has enabled automated predictive maintenance functions, analysing massive amounts of process-related data based on condition monitoring (CM). PdM stands out as the most cost-optimal maintenance type, with the potential to achieve an overall equipment effectiveness (OEE) above 90% and promising substantial returns on investment. Maintenance optimization has become a priority for industrial companies, with effective maintenance strategies capable of reducing costs significantly by addressing failures proactively (Research Paper, 2020). By leveraging Google Cloud AI Platform, our project successfully deployed a machine learning model to predict server statuses in real time, offering scalable and reliable deployment.

In summary, the literature review underscores the significant advancements in predictive maintenance brought about by the integration of ML and AI technologies. These advancements, particularly in the era of Industry 4.0, have led to a fundamental shift in how maintenance is approached, promising enhanced efficiency, reduced downtime, and overall improved operational efficacy in industrial settings. By using advanced machine learning models, web frameworks, Docker, and Google Cloud, our project showcases the future of predictive maintenance strategies that are increasingly proactive and data-driven.

Introduction

In the modern digital age, physical servers form the foundation of data centres and enterprise IT infrastructures. As businesses have expanded and digitised, the dependence on these robust servers has surged dramatically. These servers not only manage vast amounts of data but also ensure high availability, bolster security, and enhance processing speed. Consequently, the performance and reliability of these servers directly impact the efficiency of IT operations and the overall business continuity. Ensuring the optimal functioning of these servers is imperative to maintain operational efficiency and minimise disruptions. However, like all hardware, servers have a finite lifespan, but timely maintenance can maximise their operational longevity. Thus, saving companies the cost of potentially replacing their whole infrastructure due to sudden hardware failures.

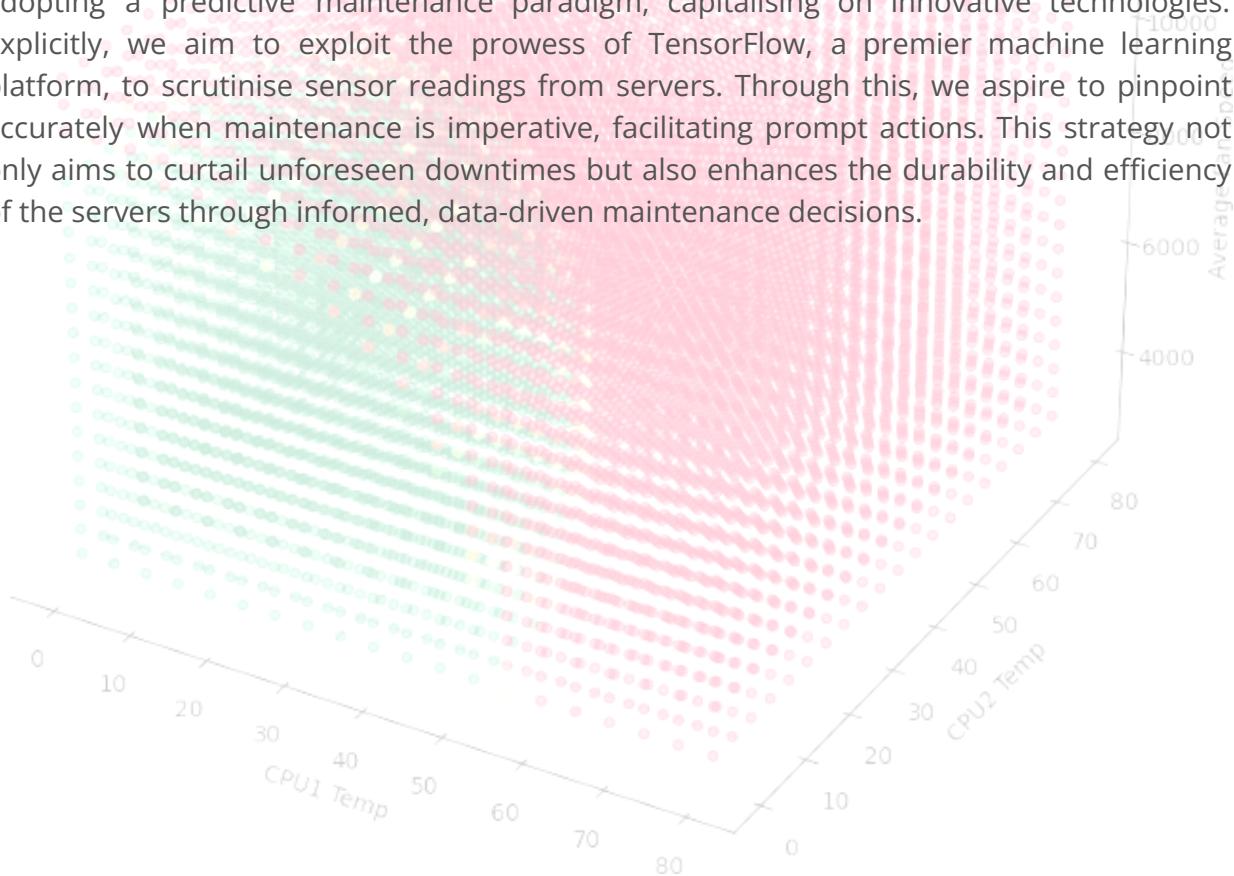
Outlined in Jack C.P. Cheng's article "Data-driven predictive maintenance planning framework for MEP components based on BIM and IoT using machine learning algorithms." Reactive maintenance, which addresses problems only post-occurrence, inherently possesses significant drawbacks. Unlike preventive or predictive maintenance strategies, it fails to proactively detect or rectify potential issues before they intensify. Sole reliance on reactive maintenance can lead to unforeseen server downtimes, as it overlooks the gradual degradation that servers experience over time. This method can often culminate in more severe failures, escalated repair expenses, and extended outages, as challenges are tackled only post-manifestation.

Preventive maintenance, though vital, falls short in precisely predicting the future health of server components. This shortfall emerges because myriad unpredictable variables can affect the performance and wear of these components. Hence, while preventive actions can alleviate potential challenges, they cannot proactively mend or adjust components based on anticipated conditions. Instituting such maintenance is not only pivotal for immediate sustenance but also for extending the overall lifespan of the servers.

Predictive maintenance, underpinned by Artificial Intelligence (AI), stands as a revolutionary advancement in the realm of server upkeep. Shown by Tyagi, V. et al. (2020), Unlike its reactive and preventive counterparts, predictive maintenance delves deep into historical and real-time data, employing sophisticated AI algorithms to discern patterns and anomalies that might be imperceptible to human analysis. By continuously monitoring server health and analysing vast datasets, AI-driven predictive maintenance can forecast potential issues long before they manifest. This not only allows IT teams to intervene

proactively, minimising disruptions, but also tailors' maintenance schedules based on actual server conditions rather than generic timelines. The integration of AI transforms the maintenance paradigm from a schedule-driven approach to a data-driven one, ensuring that servers operate at their peak efficiency while significantly reducing unforeseen downtimes. In essence, AI-augmented predictive maintenance heralds a future where server malfunctions are anticipated and mitigated, ensuring seamless and optimised IT operations.

This research project endeavours to design and craft software that integrates Machine Learning methodologies to forecast server maintenance necessities. This investigation seeks to overcome the intrinsic constraints of conventional maintenance methodologies by adopting a predictive maintenance paradigm, capitalising on innovative technologies. Explicitly, we aim to exploit the prowess of TensorFlow, a premier machine learning platform, to scrutinise sensor readings from servers. Through this, we aspire to pinpoint accurately when maintenance is imperative, facilitating prompt actions. This strategy not only aims to curtail unforeseen downtimes but also enhances the durability and efficiency of the servers through informed, data-driven maintenance decisions.



Aims and Objectives

Aim:

The primary aim of this project is to develop a predictive maintenance software solution using machine learning (ML) that can accurately predict server maintenance needs based on self attained and real-time sensor data. This project seeks to harness the power of TensorFlow to build a predictive model that can analyse key server metrics such as CPU temperature, fan speed, and error logs to forecast potential failures. The goal is to transform maintenance from a reactive or scheduled process to one driven by data, enabling businesses to optimise server longevity and operational efficiency.

Objectives:

Data Collection and Preprocessing:

Objective: To gather comprehensive sensor data from enterprise servers and preprocess it for analysis.

Justification: The quality of the predictive model depends heavily on the input data. Accurate and relevant data ensures that the model captures the patterns and anomalies necessary to predict maintenance needs effectively.

Exploratory Data Analysis:

Objective: To identify patterns and trends within the collected data that correlate with server health and failures.

Justification: Understanding the relationships within the data provides insight into which features are most predictive of failures, guiding the feature engineering and model design processes.

Model Design and Development:

Objective: To design a predictive maintenance model using TensorFlow, focusing on binary classification to determine if a server needs maintenance.

Justification: Binary classification allows for a straightforward assessment of server health, simplifying the interpretation of results and enabling quick decision-making to prevent downtime.

Model Training and Validation:

Objective: To train the predictive model using the collected data and validate its performance through testing.

Justification: Training the model on real-world data ensures it learns patterns specific to the target environment, while validation confirms the model's accuracy and ability to generalise to unseen data.

Integration and Deployment:

Objective: To deploy the predictive model in a real-time environment, incorporating it into an AI platform for live server monitoring.

Justification: Integrating the model into a real-time monitoring system allows IT teams to receive predictive maintenance alerts, enabling them to act proactively to minimise disruptions and optimise server efficiency.

Analysis of Practical Impact:

Objective: To evaluate the impact of the predictive maintenance model on server uptime, maintenance costs, and operational efficiency.

Justification: Analysing the potential and practical benefits of the model will demonstrate its value in real-world applications, supporting broader adoption of predictive maintenance strategies in the industry.

Introduction into Essential ML Tools

In the realm of Machine Learning (ML), several tools and libraries have become standard for professionals and enthusiasts alike due to their powerful features and ease of use. Here's an overview of some essential tools, each pivotal in the ML workflow:

TensorFlow

TensorFlow is an end-to-end open-source platform designed for machine learning. Developed by the Google Brain team, it has grown to be one of the most widely used ML libraries in the industry. TensorFlow excels in providing a comprehensive toolkit for researchers and developers to develop advanced ML models. One of its standout features is the ability to build and train neural networks to detect and decipher patterns and correlations, analogous to learning and reasoning used by humans. It supports a range of tasks from regression, classification, and prediction, all the way to more complex functions like natural language processing and image recognition.

TensorFlow's architecture allows for deployment across a variety of platforms (CPUs, GPUs, and even TPUs). It offers multiple abstraction levels for choosing the right one for your needs – from direct TensorFlow API commands that allow for intricate operation control to high-level Keras API which facilitates common model design patterns with ease. The flexibility and scalability of TensorFlow make it suitable not just for research and development but also for production deployment.

NumPy

NumPy is the foundational package for scientific computing in Python. It offers a powerful N-dimensional array object, sophisticated functions, tools for integrating C/C++ and Fortran code, and useful linear algebra, Fourier transform, and random number capabilities. For machine learning practitioners, NumPy is indispensable for data manipulation and preprocessing. It enables numerical operations on large data sets with speed and efficiency that native Python data structures cannot match, due to its underlying C-optimised code.

NumPy arrays form the backbone of nearly all data structures used in machine learning models, providing a much more efficient way to store and manipulate data than traditional Python lists. By facilitating operations on large arrays and matrices, NumPy serves as the bedrock upon which other libraries, including pandas and scikit-learn, are built.

Pandas

Pandas is a critical tool in the data scientist's toolkit, designed to work with structured data intuitively. It is particularly well-suited for data manipulation and analysis, offering data structures like DataFrame and Series, which are not only easy to use but also powerful for handling real-world data. pandas support a variety of data formats, allowing for easy data import, export, and manipulation.

With pandas, data scientists can perform tasks ranging from data cleaning and transformation to more sophisticated operations like data aggregation and time-series analysis. Its merging and joining capabilities are especially useful for combining datasets in complex ways, facilitating more in-depth analysis and modelling.

Matplotlib

Matplotlib is a versatile visualisation library in Python, capable of producing a wide range of static, animated, and interactive visualisations. In the context of machine learning, it is invaluable for exploratory data analysis, allowing practitioners to visualise trends, patterns, and outliers in the dataset. Through plots like histograms, scatter plots, and line charts, Matplotlib helps in understanding the data's underlying distribution, correlations, and structure.

Effective visualisation is crucial not only for exploratory analysis but also for communicating results and findings. Matplotlib provides a highly customizable interface for creating publication-quality figures and graphics that can convey complex data insights in a comprehensible and visually appealing format.

Scikit-learn

Scikit-learn is a premier library providing efficient tools for machine learning and statistical modelling including classification, regression, clustering, and dimensionality reduction. Built on NumPy, SciPy, and Matplotlib, scikit-learn offers an accessible yet versatile framework for data mining and data analysis.

Its appeal lies in its easy-to-use API and comprehensive documentation that guides users through the various algorithms it supports. With functions for fitting models, data preprocessing, cross-validation, and many more, it is designed to interoperate seamlessly with NumPy and pandas, making it a linchpin in the Python data science stack. Its

consistent interface across different types of algorithms simplifies the process of experimenting with and deploying various models, making it an ideal toolkit for both novice data scientists and seasoned practitioners alike.

Flask

Flask is a micro web framework written in Python that is lightweight and easy to use. In machine learning, Flask is often used to develop web-based APIs that can serve model predictions. Its simplicity and flexibility make it suitable for developing rapid prototypes and serving lightweight APIs for machine learning models. Flask enables machine learning practitioners to deploy their models quickly as RESTful APIs, providing easy integration with other systems.

Docker

Docker is an open platform for developing, shipping, and running applications in containers. For machine learning projects, Docker allows you to package models, scripts, and dependencies into a single portable container that can run on any platform. This ensures consistency across different environments and simplifies deployment, making it easier to share and collaborate on ML projects.

Google Cloud AI Platform

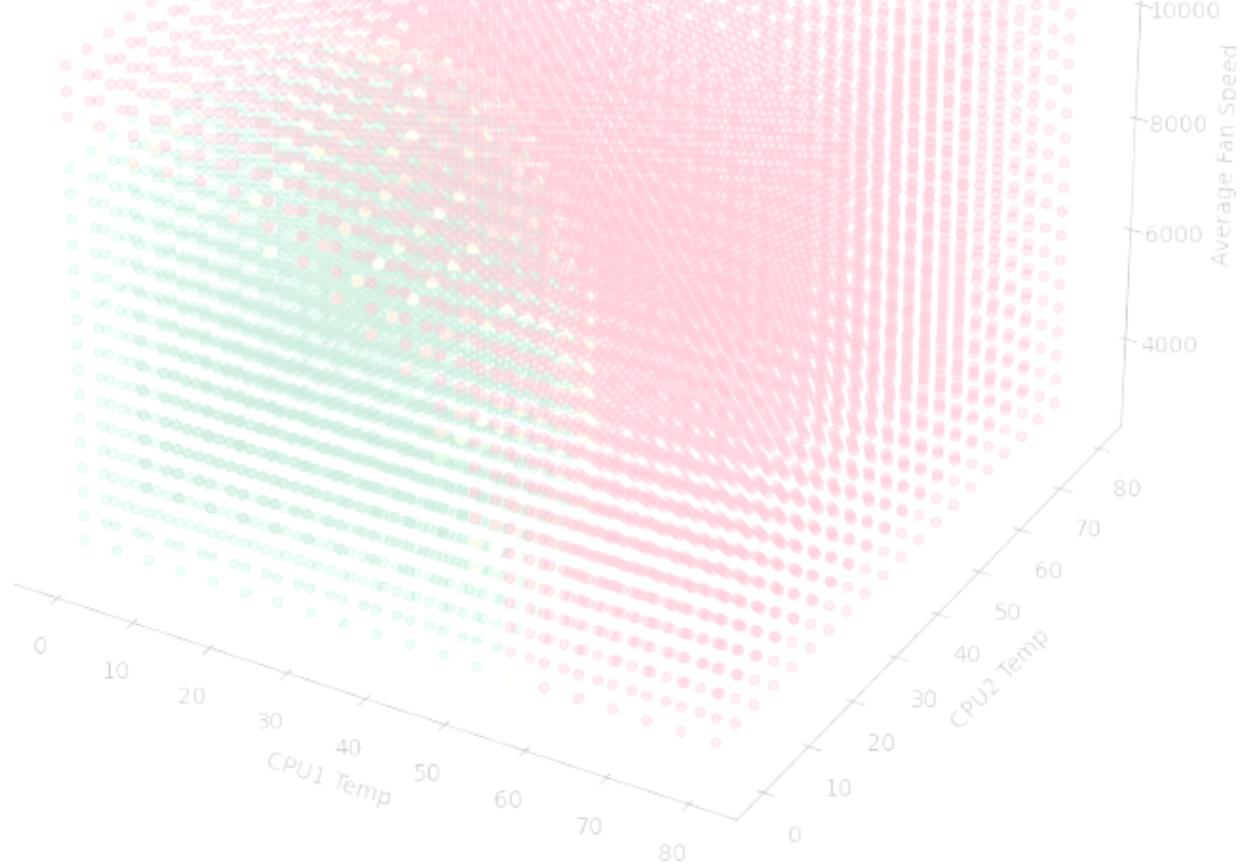
Google Cloud AI Platform is a suite of cloud services designed for deploying and managing machine learning models at scale. By integrating with Google Container Registry, it allows machine learning practitioners to deploy models in containers. This platform supports the entire ML workflow, from data preparation to model serving, and provides a scalable infrastructure to run prediction services.

GitHub

GitHub is a critical platform for version control and collaboration in machine learning projects, offering seamless integration with Jupyter Notebooks to host, manage, and share code and data. This combination enhances project transparency, reproducibility, and teamwork by centralising workflows in a single repository where changes are tracked, facilitating easy revisions and collaboration without data loss or overwrite concerns. GitHub

not only streamlines code and project management through features like pull requests and issue tracking but also fosters an open-source community where data scientists and developers share knowledge and advancements. Incorporating GitHub in machine learning workflows using Jupyter Notebooks ensures that projects are not just technically robust but also well-documented and accessible, promoting shared learning and innovation in the ML community.

By deeply understanding each of these tools, machine learning practitioners can effectively tackle tasks ranging from data handling and processing to model development and evaluation, leveraging their unique capabilities to enhance both the efficiency and quality of their work.

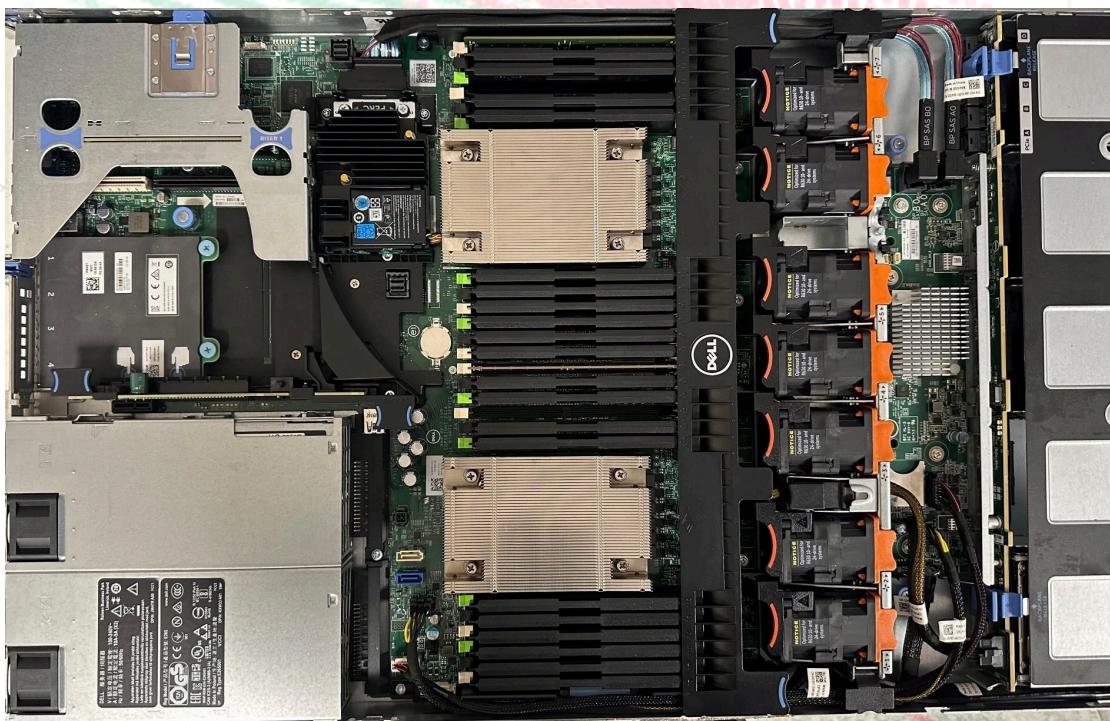


Servers

Video: <https://youtu.be/t-vj5j1CkvU>



Servers are foundational to the modern world of Information Technology (IT) and beyond, playing a critical role in data processing, storage, and management across various sectors including finance, healthcare, education, and e-commerce. They are the backbone of the internet, hosting the applications and services we use daily, from cloud computing and online banking to streaming services and social media platforms. Servers ensure that data is accessible, secure, and efficiently managed, enabling businesses and organisations to operate seamlessly, analyse big data for insights, and provide the digital services that have become integral to our daily lives.



In the context of maintaining these vital systems, the use of sensors within servers becomes critical. Sensors monitor various parameters, such as CPU temperature and fan speed, providing real-time data that is essential for maintaining operational efficiency and preventing downtime. For instance, overheating can significantly impair a server's performance and, in severe cases, lead to hardware damage. By monitoring CPU temperature and fan speeds, IT professionals can intervene early to mitigate risks, such as by improving ventilation or performing maintenance tasks.

The integration of TensorFlow-based AI algorithms into predictive maintenance methodologies represents a significant advancement in optimising the operational longevity and efficiency of physical servers in data centres and enterprise IT infrastructures. TensorFlow's ability to analyse complex data sets enables the development of models that can predict potential failures or identify inefficiencies in server operations before they become critical issues. For example, by analysing trends in temperature data and fan speed, TensorFlow can predict when a server is likely to overheat or when a fan is failing, allowing for pre-emptive maintenance actions that can avoid costly downtime and extend the server's lifespan.

Furthermore, TensorFlow's machine learning capabilities can optimise workload distribution based on the thermal behaviour of individual CPUs within a server. This ensures that no single CPU is overburdened, reducing the risk of overheating, and improving overall system efficiency. By leveraging such AI-driven insights, Data Centres can significantly enhance their predictive maintenance strategies, leading to more reliable, efficient, and cost-effective operations.

The criticality of servers in IT and the importance of sensor data for their maintenance underscores the value of integrating advanced AI algorithms like those offered by TensorFlow. This integration not only enhances the ability to maintain and optimise server operations but also represents a forward-looking approach to managing the increasingly complex and crucial IT infrastructures that support our digital world.

Central Processing Unit (CPU)

The CPU is essentially the brain of the computer, handling millions of processes per second. High temperatures can degrade its performance over time or cause immediate throttling, where the CPU reduces its speed to prevent overheating. This throttling can lead to slower system performance and, in severe cases, system crashes or hardware damage. Monitoring CPU temperature helps in early detection of potential overheating issues, allowing for timely intervention like cleaning dust from the system, improving ventilation, or replacing the thermal paste.



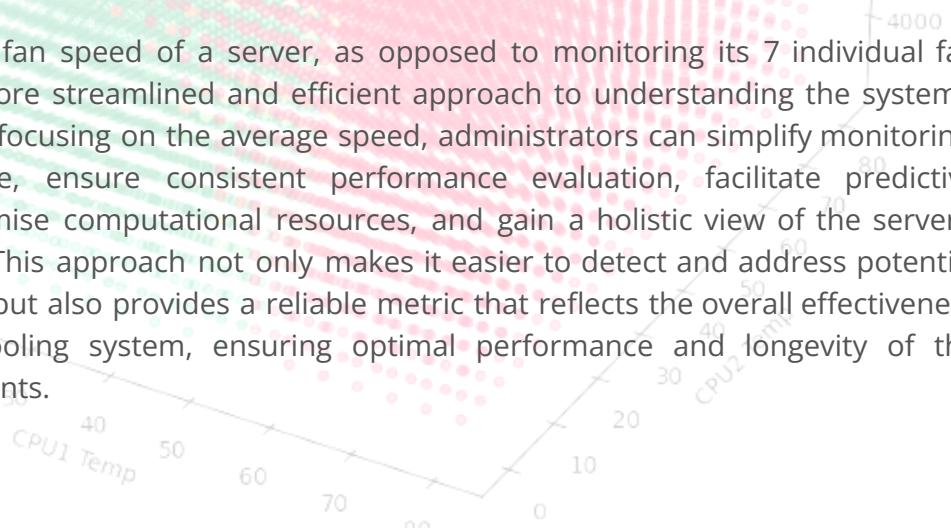
Monitoring each individual CPU's temperature within a server, rather than relying on an average temperature, provides more granular insight into the thermal status of each processing unit, enabling precise identification of localised overheating issues that an average temperature might mask. This detailed monitoring is crucial because even if one CPU overheats while others remain cool, the average temperature could appear normal, potentially overlooking critical hotspots that can lead to CPU throttling or failure. Individual temperature readings allow for targeted cooling adjustments and more effective troubleshooting, ensuring that all CPUs operate within their thermal thresholds for optimal performance and reliability. Additionally, understanding the specific thermal behaviour of each CPU can inform better workload distribution and cooling strategies, enhancing the overall efficiency and longevity of the server.

Fans

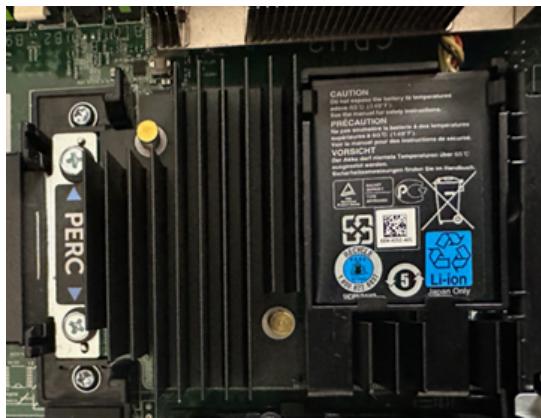
Fans play a crucial role in maintaining optimal operating temperatures by dissipating heat away from critical components like the CPU and GPU. The fan speed adjusts according to the system's cooling needs; higher temperatures typically demand higher fan speeds. Monitoring fan speeds can provide insights into the overall thermal status of the system. Unusually high fan speeds might indicate excessive heat or inefficient cooling, while unusually low speeds could signal fan malfunctions or obstructions that impede airflow.



Using the average fan speed of a server, as opposed to monitoring its 7 individual fan speeds, offers a more streamlined and efficient approach to understanding the system's thermal health. By focusing on the average speed, administrators can simplify monitoring, reduce data noise, ensure consistent performance evaluation, facilitate predictive maintenance, optimise computational resources, and gain a holistic view of the server's cooling efficiency. This approach not only makes it easier to detect and address potential issues proactively, but also provides a reliable metric that reflects the overall effectiveness of the server's cooling system, ensuring optimal performance and longevity of the hardware components.



RAID Controller (Redundant Array of Independent Disks)



The Redundant Array of Independent Disks (RAID) card is crucial for data redundancy and performance enhancement. It allows multiple hard drives to work together, improving the overall system's fault tolerance and data integrity. In a data collection environment, where the loss or corruption of data can have significant consequences, a RAID setup ensures that data is mirrored across multiple drives. This means that if one drive fails, the system can continue to operate without data loss, and the failed drive can

be replaced without downtime. Additionally, RAID can be configured to enhance the read/write speed, essential for the high-speed data transactions typical in server operations.

RAM (Random Access Memory)



Random Access Memory (RAM) is vital for the temporary storage of data that the server's processor needs to access quickly. High-capacity, high-speed RAM is essential for efficient data processing, allowing for faster retrieval and manipulation of data. This is particularly important in data collection scenarios where the server must handle large datasets or run multiple applications simultaneously. Sufficient RAM ensures that these operations can be performed smoothly, without lag or bottlenecks, significantly affecting the server's ability to collect, process, and analyse data efficiently.

Network Interface Controller



A NIC is fundamental for establishing and managing the server's connection to a network. In the context of data collection, a high-performance NIC ensures that the server can handle high volumes of data ingress and egress without network bottlenecks. It is responsible for the fast and reliable transmission of data between the server and other networked devices or internet-based resources, making it a key component for servers that rely on network-intensive applications or

need to transmit collected data to remote storage or analysis services.

Hard Disk Drive

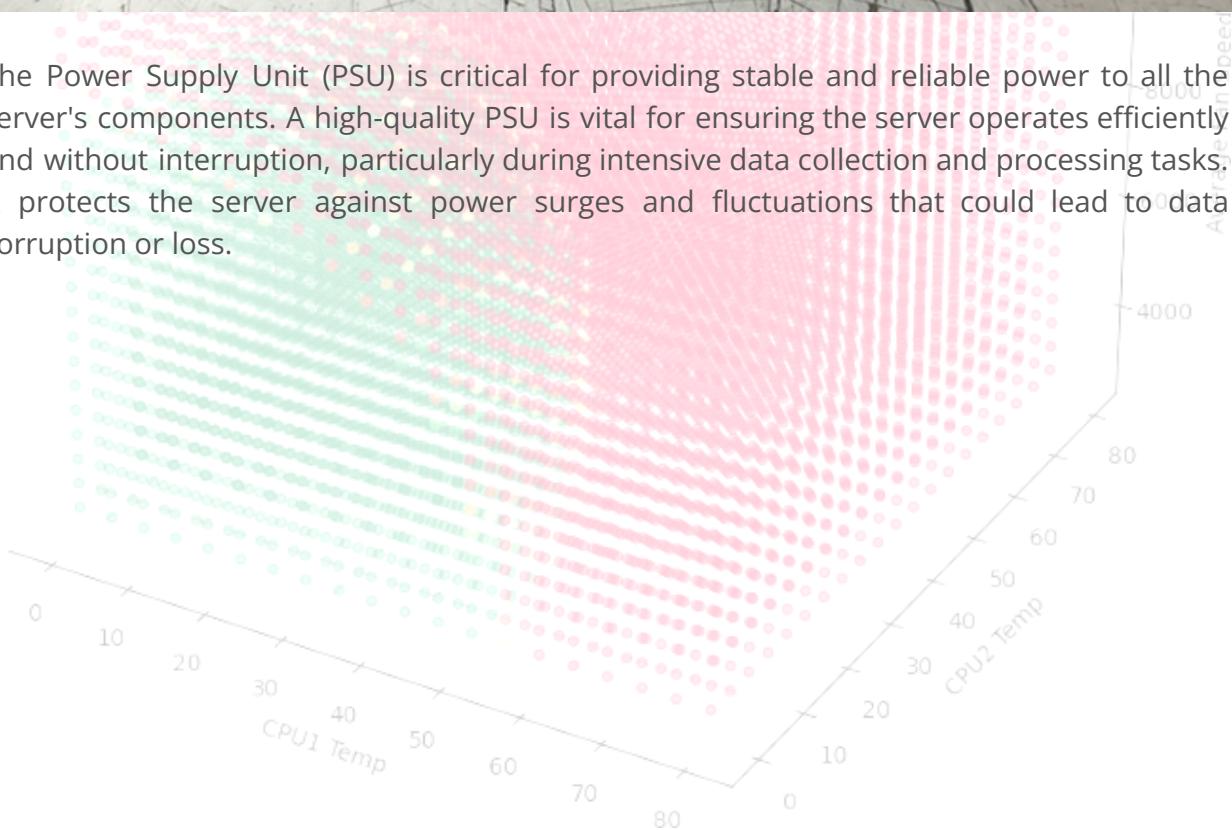
The hard drive serves as the primary storage device for the server, holding the operating system, applications, and, most importantly, the collected data. In a data collection environment, the choice of hard drive impacts the server's storage capacity, speed, and reliability. Solid-state drives (SSDs) offer faster data access speeds than traditional hard disk drives (HDDs), making them preferable for situations where speed is crucial. However, HDDs may still be used for long-term storage of large volumes of data where speed is less critical. The selection between SSDs and HDDs (or a combination of both) depends on the specific needs of the data collection task, including considerations of speed, capacity, and cost.



PSU (Power Supply Unit)



The Power Supply Unit (PSU) is critical for providing stable and reliable power to all the server's components. A high-quality PSU is vital for ensuring the server operates efficiently and without interruption, particularly during intensive data collection and processing tasks. It protects the server against power surges and fluctuations that could lead to data corruption or loss.



Overall

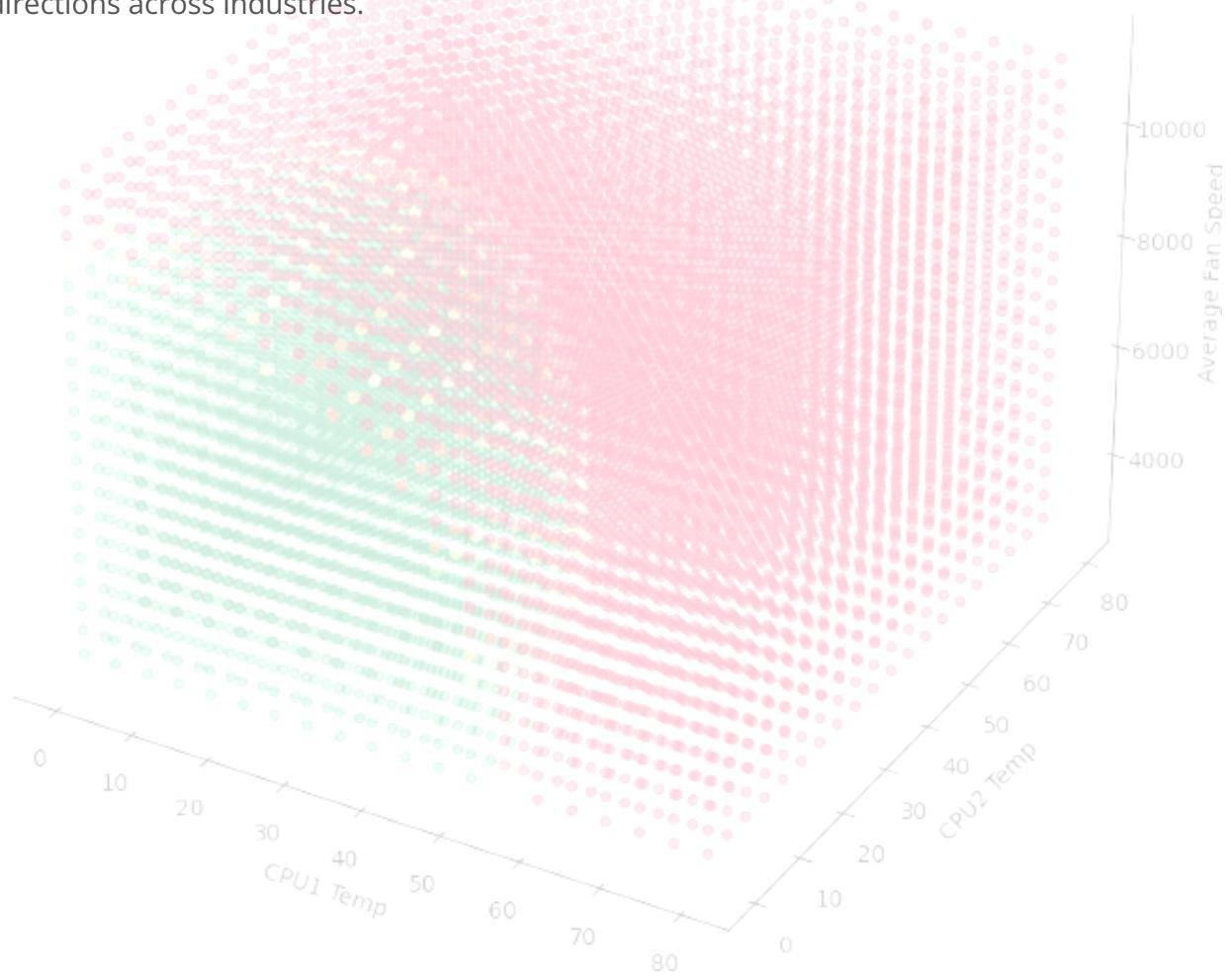
This study, focusing on the use of TensorFlow-based AI algorithms for predictive maintenance, underscores the immense potential of integrating advanced computational techniques within server environments. This integration is crucial not just for servers like the Dell PowerEdge R630, which is used for the sensor data collection of this study, but also across various models employed by different IT companies such as the diverse configurations of the Dell PowerEdge R930, HP ProLiant Gen9 DL580, and Dell PowerEdge R730.



Moreover, the comprehensive justification for data collection within server environments, as detailed through the exploration of critical components such as the RAID card, RAM, NIC, PSU, and hard drives, underscores the intricate balance between hardware reliability, efficiency, and the advanced computational requirements of modern IT infrastructures. Each component, from ensuring data redundancy and enhancing processing speed to facilitating robust network connections and providing dependable storage solutions, plays a pivotal role in the overarching goal of optimising server operations for the demanding tasks of data collection and analysis.



This exploration not only highlights the necessity of each hardware component in maintaining the operational integrity and performance of servers but also illuminates the potential for integration with cutting-edge technologies like TensorFlow-based AI algorithms. Such integration promises to revolutionise predictive maintenance methodologies, further enhancing the resilience, efficiency, and longevity of servers in data centres and enterprise IT environments. By leveraging these technologies and insights, organisations can anticipate and preemptively address potential issues, ensuring the continuous, reliable operation of their servers. This proactive approach to server maintenance and optimization is essential in an era where data is not just an asset but the backbone of operational intelligence, driving decisions, innovations, and strategic directions across industries.

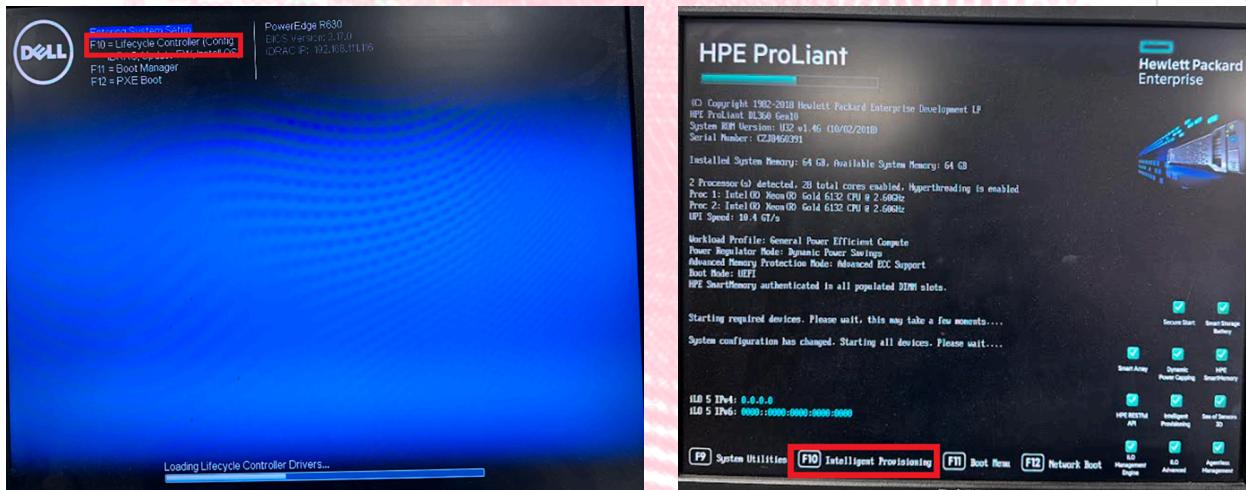


Data Collection

Video: <https://youtu.be/bEuZfevDopg>

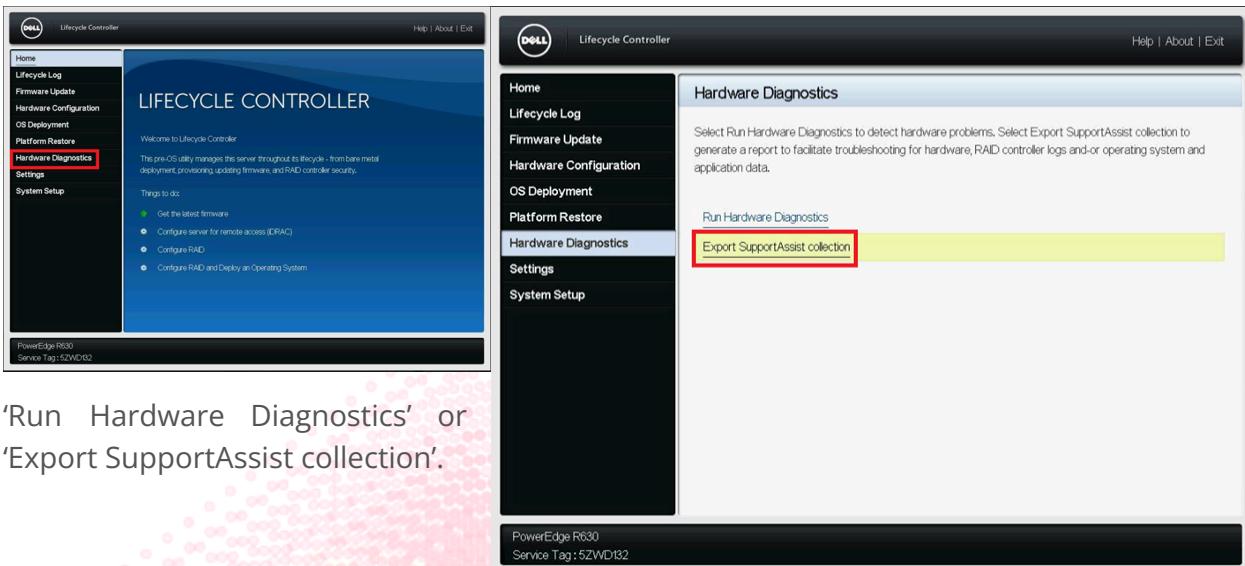
For this particular case, we will be collecting the data (CPU Temperature and fan speed) from the onboard Dell application known as the 'Lifecycle Controller' which is delivered as part of integrated Dell Remote Access Controller (iDRAC) out-of-band solution and embedded Unified Extensible Firmware Interface (UEFI) applications in the latest Dell servers.

When booting up a 13th Generation Dell Server and the USB Ports have initialised, you will be presented with 4 options. For the purpose of data collection, I will only mention the option that is of concern which is the F10 Option



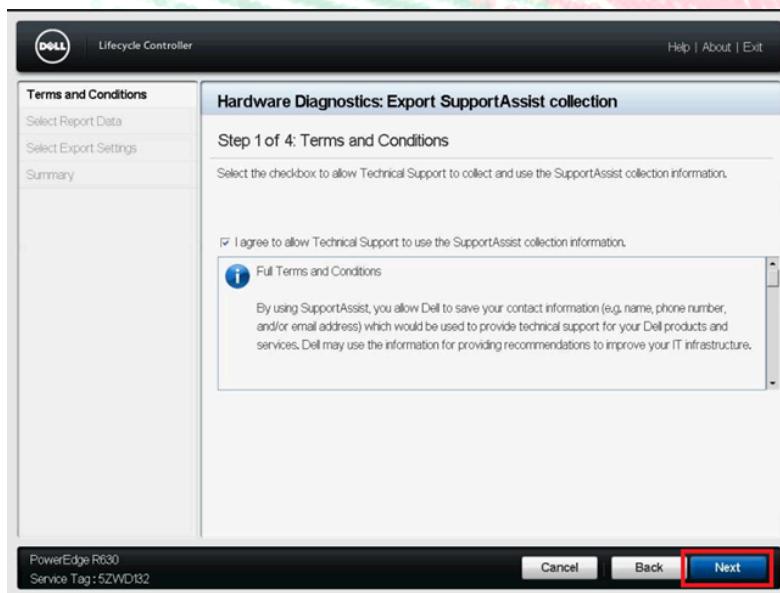
Before we continue any further, it must be mentioned that what is to be shown can also be done with many other brands of server such as IBM and HP with their own and respective management firmware, for instance, the HP Model variant would still list you with 4 options whilst booting and again you will enter F10 to enter its 'Intelligent Provisioning' firmware to then view and collect its sensor data

Going back to the Dell R630 we were looking at initially, after it has booted successfully into the Lifecycle Controller, we will be presented with a host of utilities to enter and in this case, we shall enter the 'Hardware Diagnostics' tab which will allow us to do one of two things:

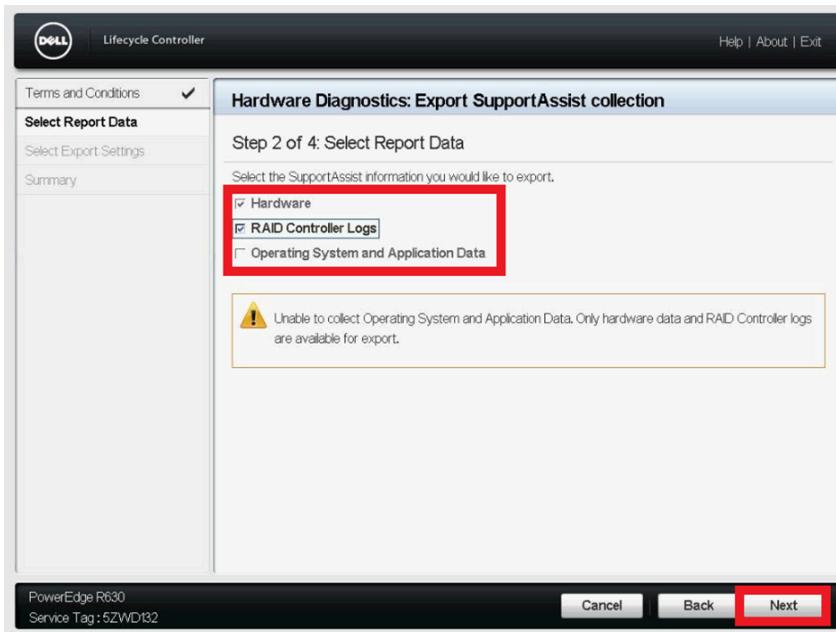


'Run Hardware Diagnostics' or
'Export SupportAssist collection'.

The 1st option will lead us to the preventative/reactive form of maintenance. This test can vary from taking between 20 minutes to 2 hours depending on the issue which can mean a long and undesirable downtime of a system but it does give a detailed breakdown of all tests and results when the server is put under stress or not. We want to export the initial/running diagnosis of all parts of the server where a sensor is located which not only gives us valuable information but also does not take a long time to do (usually around 4 minutes) and with this data, we will use it to build the predictive form of maintenance with the use of TensorFlow and AI.

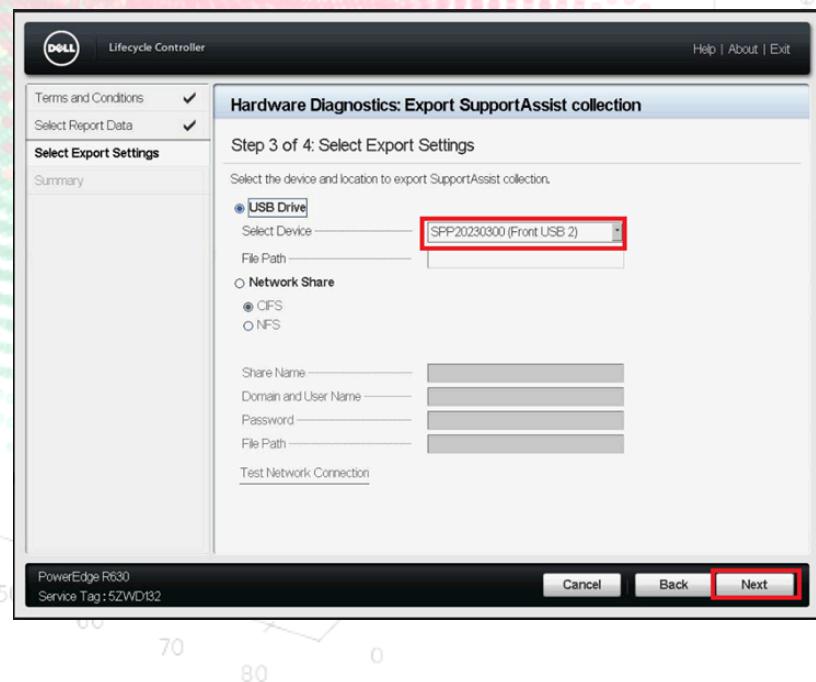


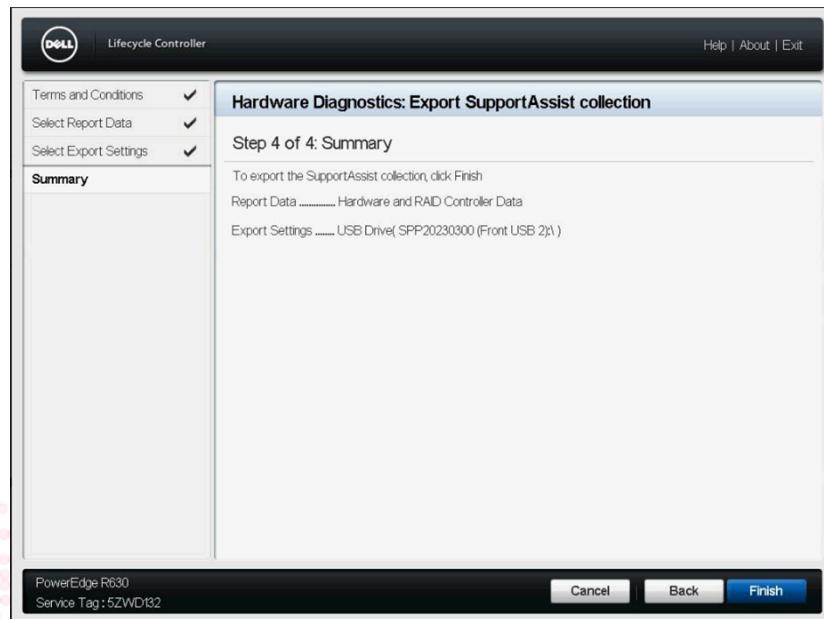
Click "I Agree" to continue ->



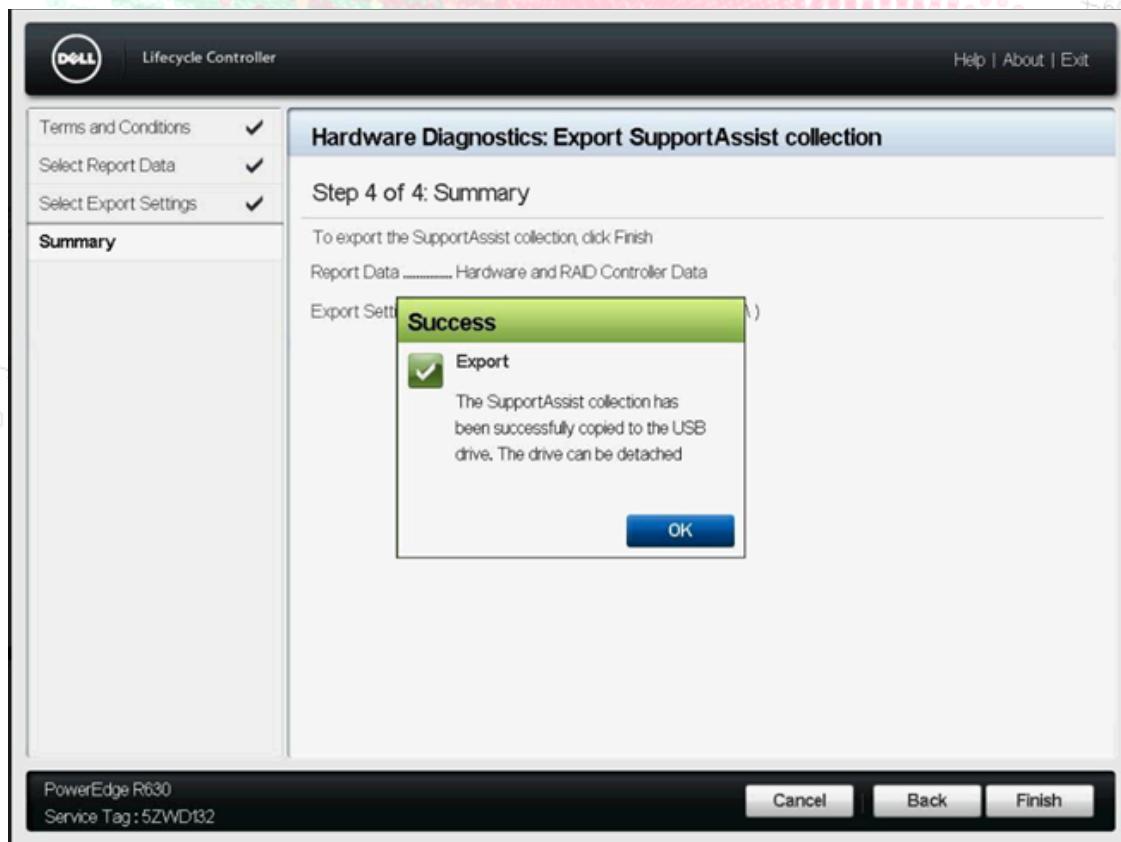
This screen dictates where the gathered information should be stored, I chose using a USB for simplicity but via the NFC or CIFS protocol, the information can be directly stored on your computer, or any other virtual drive located.

All boxes that can be checked must be checked to ensure we collect the maximum amount of sensor data possible then click next to continue again. ->





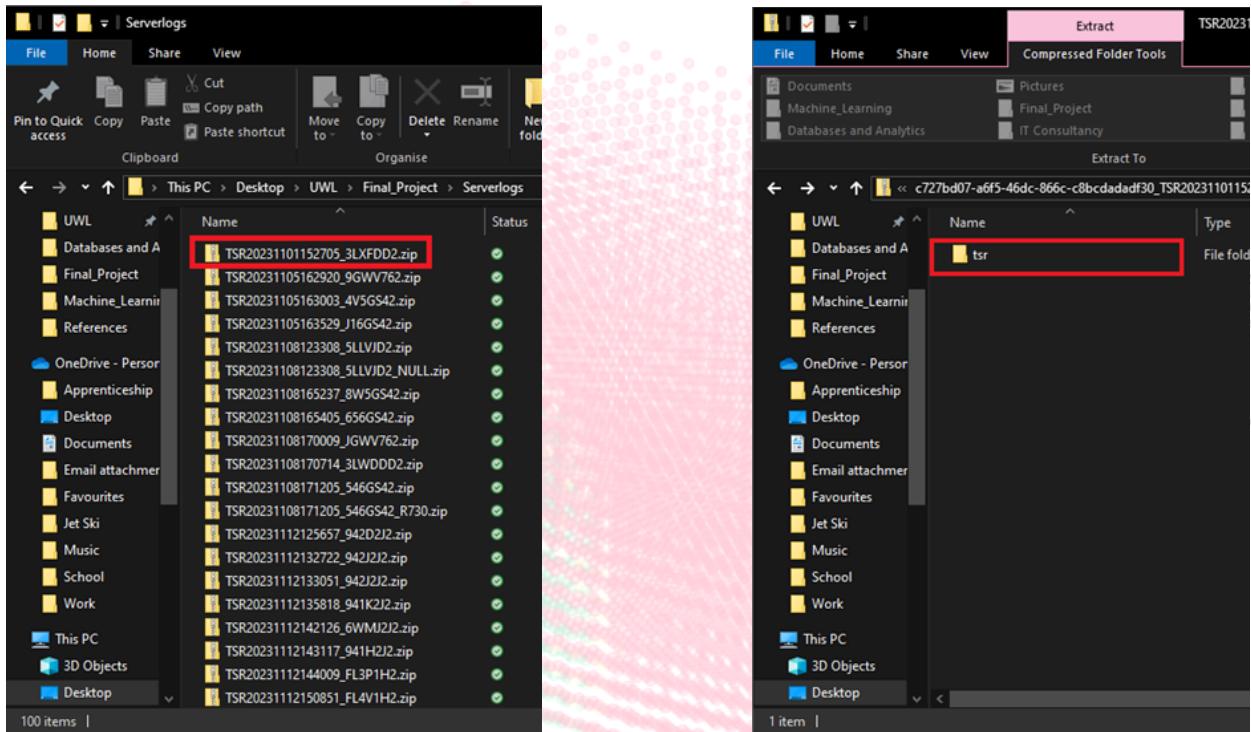
Click finish on the final confirmation screen to then complete this process for 1 server, this whole process must be completed for each server (Data Point)



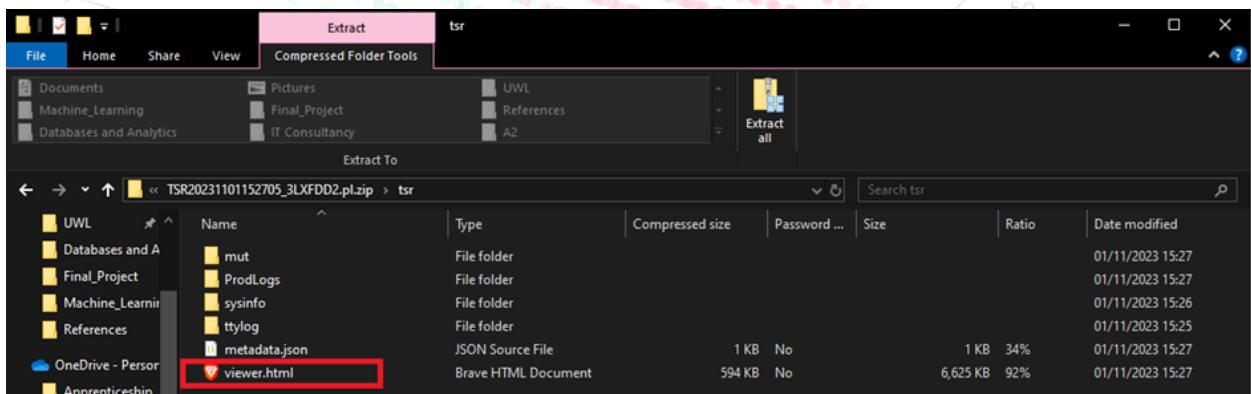
Data Organisation

Video: <https://youtu.be/yHOxb2ICOXU>

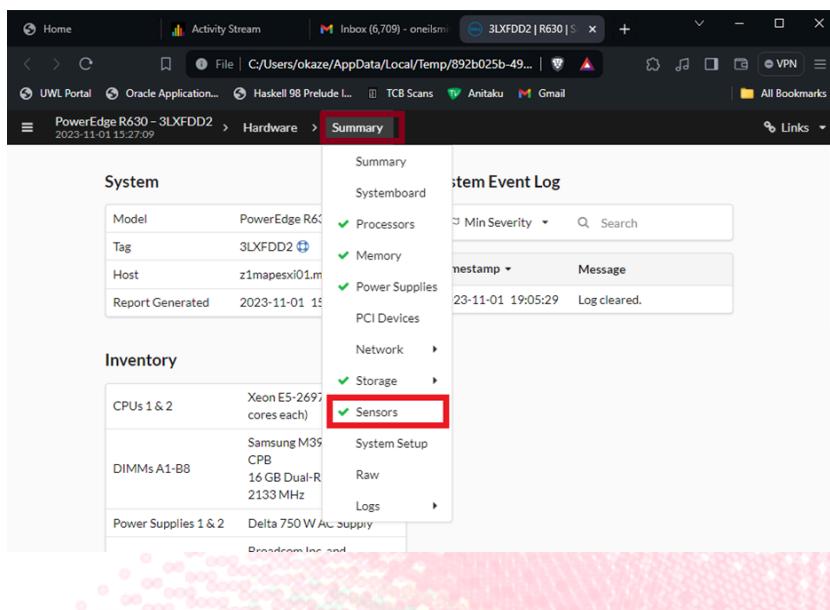
After data has been exported from the physical server, it appears in the form of a zip file that when extracted, provides a html viewer of all collected information



Once extracted you can now begin to view all information provided by the server.

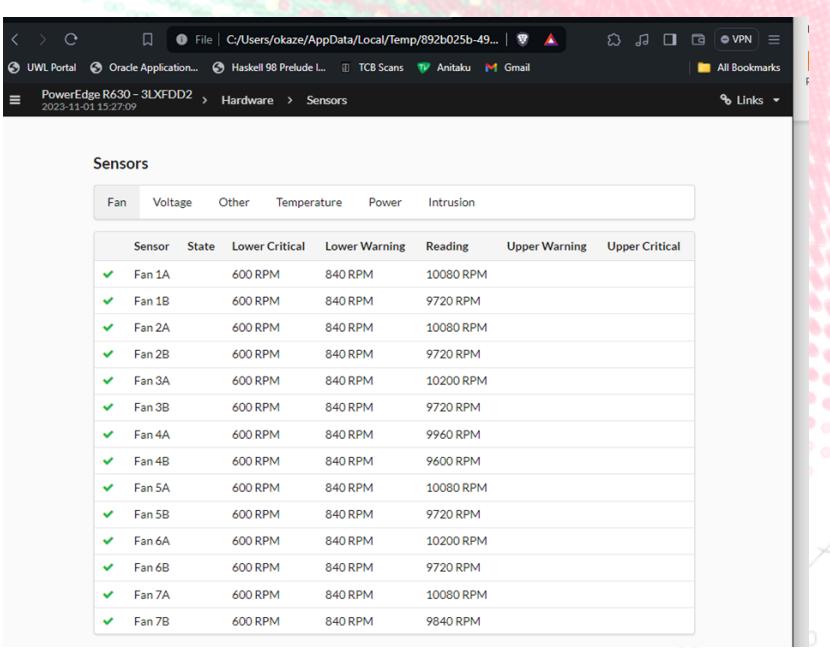


Clicking on "viewer.html" will then take you to the next screen which contains all the data it a neatly contained fashion



The screenshot shows the 'Summary' tab of the PowerEdge R630 hardware interface. It displays the following information:

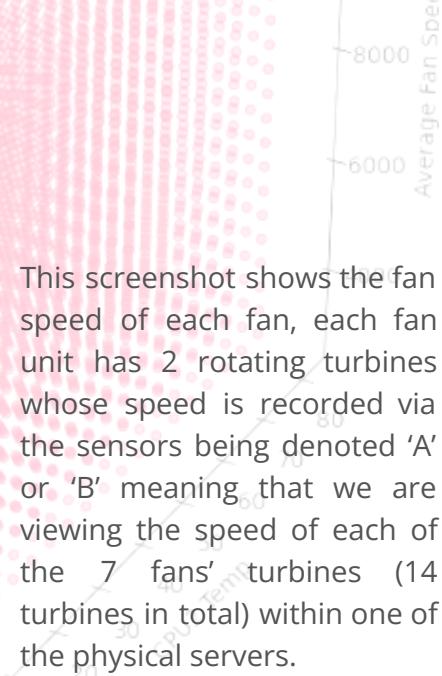
- System:**
 - Model: PowerEdge R630
 - Tag: 3LXFDD2
 - Host: z1mapesx01.m
 - Report Generated: 2023-11-01 15:27:09
- Inventory:**
 - CPU: Xeon E5-2697 v4 (24 cores each)
 - DIMM: Samsung M39 CPB 16 GB Dual-R Rank 2133 MHz
 - Power Supply: Delta 750 W AC-supply
- Sensors:** A dropdown menu under the 'Hardware' tab is open, with 'Sensors' highlighted.



The screenshot shows the 'Sensors' page of the PowerEdge R630 interface. It displays the following data:

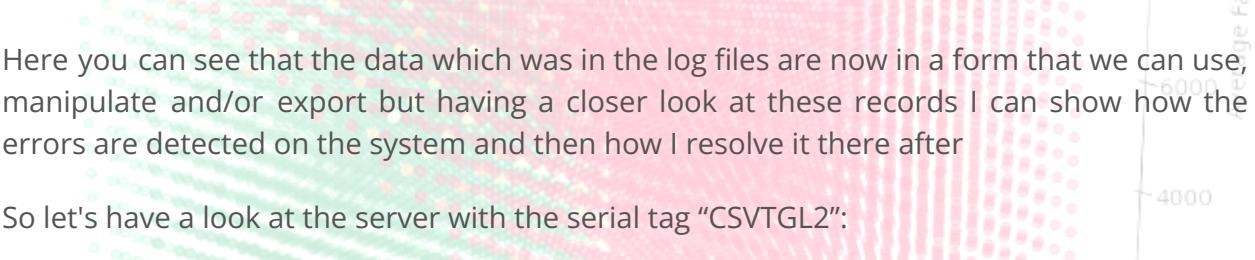
Fan	Voltage	Other	Temperature	Power	Intrusion
Fan 1A	600 RPM	840 RPM	10080 RPM		
Fan 1B	600 RPM	840 RPM	9720 RPM		
Fan 2A	600 RPM	840 RPM	10080 RPM		
Fan 2B	600 RPM	840 RPM	9720 RPM		
Fan 3A	600 RPM	840 RPM	10200 RPM		
Fan 3B	600 RPM	840 RPM	9720 RPM		
Fan 4A	600 RPM	840 RPM	9960 RPM		
Fan 4B	600 RPM	840 RPM	9600 RPM		
Fan 5A	600 RPM	840 RPM	10080 RPM		
Fan 5B	600 RPM	840 RPM	9720 RPM		
Fan 6A	600 RPM	840 RPM	10200 RPM		
Fan 6B	600 RPM	840 RPM	9720 RPM		
Fan 7A	600 RPM	840 RPM	10080 RPM		
Fan 7B	600 RPM	840 RPM	9840 RPM		

Once arrived at this screen, it can be clearly seen that a whole host of information from the server has been made available to be viewed, from what OS the server is running to which hardware is keeping the server operational. But what we want to see is the data returned by the sensors embedded within the server hence we click on sensors after the summary tab has been clicked



This screenshot shows the fan speed of each fan, each fan unit has 2 rotating turbines whose speed is recorded via the sensors being denoted 'A' or 'B' meaning that we are viewing the speed of each of the 7 fans' turbines (14 turbines in total) within one of the physical servers.

Once all necessary data has been retrieved from the servers, I then keep note of them by manually inputting the data into an excel workbook as shown below.



A screenshot of an Excel spreadsheet titled "PM_SensorData.xlsx". The spreadsheet contains data from multiple log files, specifically focusing on temperature and fan speed measurements. The columns include "Server Tag", "CPU1 TEMP", "CPU2 TEMP", "CPU3 TEMP", "CPU4 TEMP", "CPU5 TEMP", "FAN 1A", "FAN 1B", "FAN 2A", "FAN 2B", "FAN 3A", "FAN 3B", "FAN 4A", "FAN 4B", "FAN 5A", "FAN 5B", "FAN 6A", "FAN 6B", "FAN 7A", "FAN 7B", "AVG FAN", "Working?", "E-CODE", and "SOLUTION". The data shows various temperatures in Celsius and fan speeds in RPM across different server components. A red box highlights the "Tag" column, which contains the serial tag "CSVTGL2".

Here you can see that the data which was in the log files are now in a form that we can use, manipulate and/or export but having a closer look at these records I can show how the errors are detected on the system and then how I resolve it there after

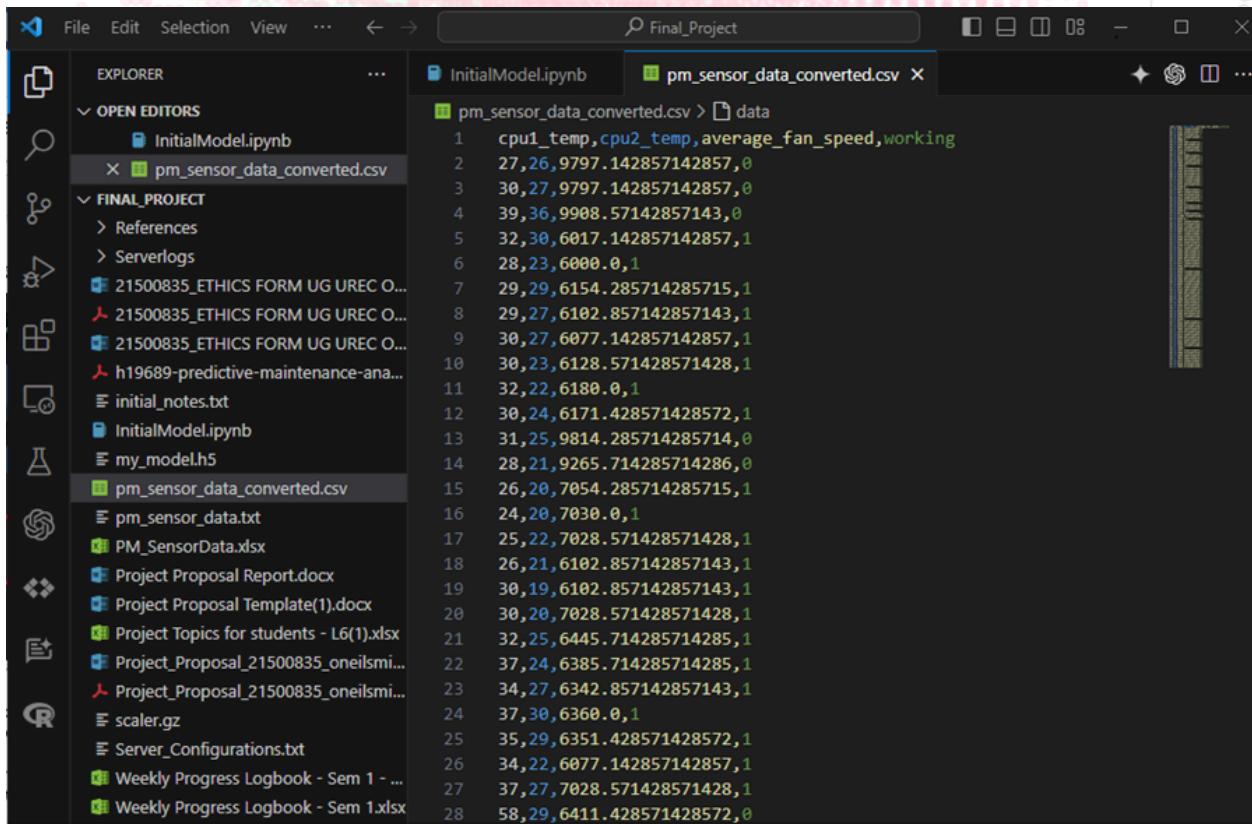
So let's have a look at the server with the serial tag "CSVTGL2":



A screenshot of the Dell SupportAssist Collection Viewer interface. The top navigation bar includes links for Home, Activity Stream, Inbox (6,711), 3LXFDD2 | R, CSVTGL2 (R), CSWTGL2 (R), and a plus sign. Below the navigation bar, there are tabs for Overview, Board, CPU, Memory, Power, PCI, Network, Storage, and Sensors. The "Config" tab is currently selected. On the left, there is a "SupportAssist Collection Viewer" sidebar with sections for System, Inventory, and Sensors. The "System" section displays the server model as "PowerEdge R630", the tag as "CSVTGL2" (which is highlighted with a red box), the host as "r1nu0.ac3sadc01.bc.jsplc.net", and the report date as "2023-12-17 14:10:00". The "Inventory" section shows two CPUs (1 & 2) as Xeon E5-2680 v4 (14 cores each) and DIMMs A1-B1 as Hynix Semiconductor. On the right, there is a "Historic SEL Entries" section showing two entries: one from 2023-12-17 at 17:40:52 stating "The chassis is closed while the power is off.", and another from 2023-12-17 at 19:46:47 stating "The chassis is open while the power is off." Both of these entries are also highlighted with a red box.

Here we can see that this server had something which is also known as an intrusion error, it detected that the chassis was not entirely closed between the moments of the server booting up or turning off. Because of this, the fans were spinning faster than necessary. This was resolved by changing the lid which was originally with this server to another that fitted appropriately, thus returning the fan speeds to an optimal speed. Looking further into this, if the fan speeds were allowed to keep spinning at this rate whilst the rest of the server is performing normally, the fans would've loosened themselves out or burnt out when the server is under peak performance conditions.

After all the data needed has been added to the excel workbook, I then use ChatGPT to convert that same excel workbook into a CSV(comma separated values) file which can then be used for a machine learning model which is developed within a Jupyter notebook format (python/.ipynb)



```

InitialModel.ipynb pm_sensor_dataConverted.csv
cpu1_temp,cpu2_temp,average_fan_speed,working
27,26,9797.142857142857,0
38,27,9797.142857142857,0
39,36,9908.57142857143,0
32,30,6017.142857142857,1
28,23,6000.0,1
29,29,6154.285714285715,1
29,27,6102.857142857143,1
30,27,6077.142857142857,1
30,23,6128.571428571428,1
32,22,6180.0,1
38,24,6171.428571428572,1
31,25,9814.285714285714,0
28,21,9265.714285714286,0
26,20,7054.285714285715,1
24,20,7030.0,1
25,22,7028.571428571428,1
26,21,6102.857142857143,1
30,19,6102.857142857143,1
30,20,7028.571428571428,1
32,25,6445.714285714285,1
37,24,6385.714285714285,1
34,27,6342.857142857143,1
37,30,6360.0,1
35,29,6351.428571428572,1
34,22,6077.142857142857,1
37,27,7028.571428571428,1
58,29,6411.428571428572,0

```

Supervised Learning

Supervised learning is a type of machine learning where models are trained using labelled data—data that includes an input (feature) and an output (label). This method leverages historical data that has been correctly annotated to learn the relationships and patterns between input features and the target output. In supervised learning, the model iteratively makes predictions on the data and is corrected by the teacher (the correct outputs), refining its algorithms until it achieves a desirable level of accuracy.

In the context of this project, supervised learning has been pivotal in developing a predictive maintenance system for servers, such as those depicted in the Dell PowerEdge R630. The sensor data collected includes critical parameters like CPU1 temperature, CPU2 temperature, average fan speed, and a categorical label indicating whether the server is functioning correctly ("working") or not ("not working") identified visually. These labels are crucial for training our model to recognize the conditions under which a server may fail, thus enabling it to predict such failures before they occur.

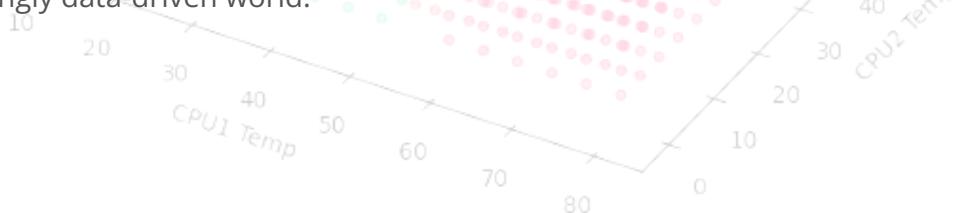
The main objective of employing supervised learning in this project is to minimise the occurrence of false positives (predicting a failure when the server is fine) and false negatives (failing to predict an impending failure). False positives can lead to unnecessary checks and maintenance, wasting resources and time, whereas false negatives can be even more detrimental, potentially leading to unexpected downtime and the associated costs and disruptions. By training the model with accurately labelled data, the predictive algorithm learns to discern subtle patterns that precede failures, which might be overlooked by traditional monitoring systems.

The model's efficacy can be continually validated through a confusion matrix—a tool that helps visualise the performance of an algorithm. Each entry in the matrix allows us to understand the number of correct and incorrect predictions made by the model, categorised by type (true positives, false positives, true negatives, and false negatives). This tool is essential for tuning the model to reduce errors and enhance its predictive accuracy, thereby supporting the proactive maintenance strategy that is central to the project's goals.

The reason for taking this approach came from Dr Massoud Zolghani's week 4 lesson "Supervised Classification and K-NN Performance Metrics" in the module of "Machine Learning"



By leveraging TensorFlow within a supervised learning framework, this project harnesses the power of AI to enhance server maintenance protocols. Not only does this approach reduce the likelihood of unexpected server failures, but it also optimises maintenance schedules, ensuring that interventions are made precisely when needed, thus extending the hardware's lifespan and maintaining the operational integrity of IT infrastructures in an increasingly data-driven world.

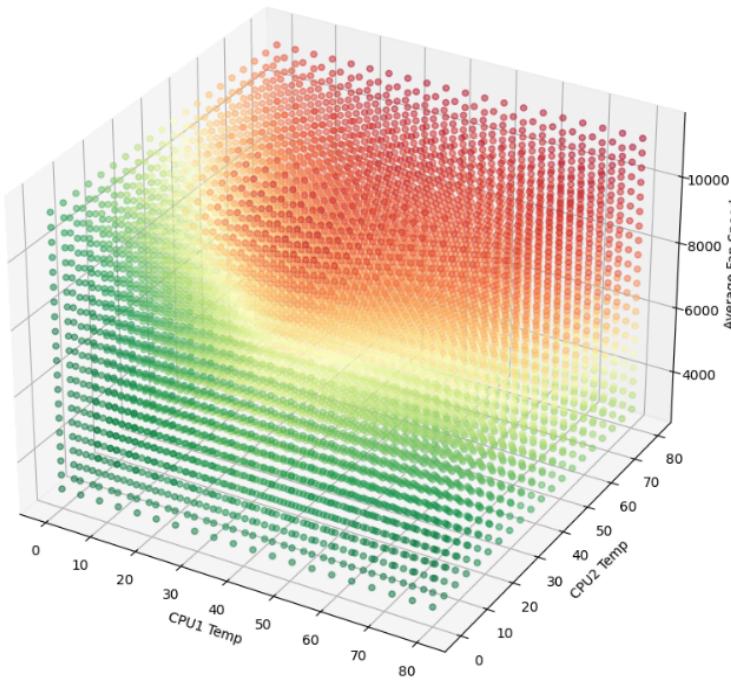


Python Script

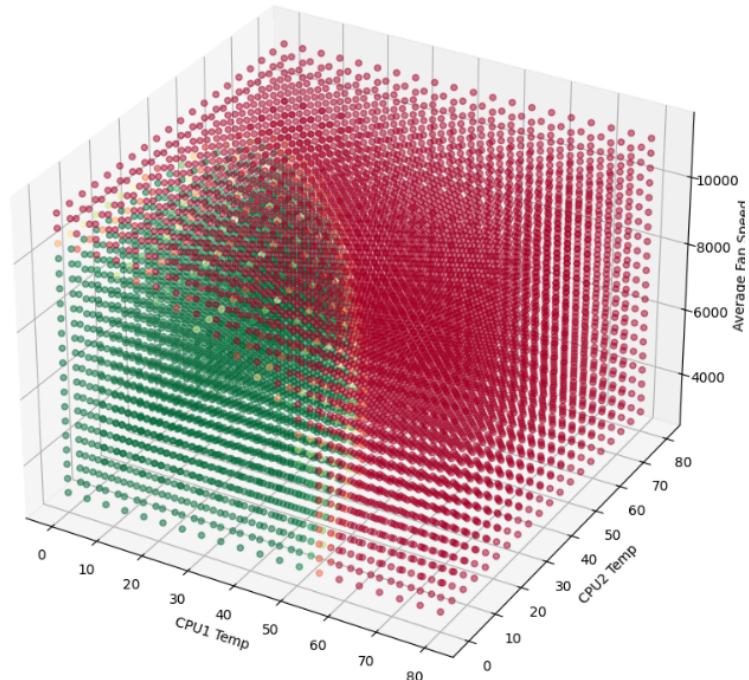
This section of the report outlines the development and functionality of the predictive maintenance models, termed as "Initial_Model" and "Final_Model". Both models are implemented in Python using Jupyter notebooks, leveraging TensorFlow for building and training machine learning models based on supervised learning techniques. These models analyse server sensor data—specifically CPU temperatures, fan speeds, and operational status—to predict maintenance needs.

With the next two sub-chapters,

3D Decision Boundary Visualization Covering All Data Points



3D Decision Boundary Visualization Covering All Data Points





Initial_Model

