

# Evaluating Machine Learning Algorithms for Fault Detection in Solar Cells: Performance and Limitations

Asst.Lec. Zaher Fadhil  
Raham<sup>1</sup>

[Zahir.raham@gmail.com](mailto:Zahir.raham@gmail.com)

Lec. Rafat K. Oubida<sup>2</sup>

[raafatoubida@nahrainuniv.edu.iq](mailto:raafatoubida@nahrainuniv.edu.iq)

Prof.Dr. Abdul Monem S. Rahma<sup>3</sup>

[monem.rahma@muc.edu.iq](mailto:monem.rahma@muc.edu.iq)

**Abstract:** The extensive distribution, the penetration and popularity of PV have turned fault diagnostics from a general inspection into an operation necessity. With a larger PV system that becomes ever more complex, undetected fault can cause cumulative power failure with lower system performance and hence cost, especially in situations where real-time monitoring is a prerequisite. Whilst infrared thermography, electroluminescence (EL) imaging, and current–voltage (I–V) curve analysis are common diagnostic methods, their automated implementation in industrial-scale PV plants remains limited by practical and computational bottlenecks. We propose a multimodal fault diagnosis framework that incorporates EL image analysis with infrared thermal measurements and electrical I–V data collected by operating PV systems. Instead of depending purely on a single diagnostic method, this framework studies the performance metrics across the learning paradigms for identical experimental conditions. Based on the same dataset, a convolutional neural network (CNN) is applied and compared with established machine learning models, i.e. support vector machines (SVM) and extreme gradient boosting (XGBoost). In experimental testing, we use 2,624 real EL images. The results are as follows: CNN models have the best level of classification at 95.2%, whereas inference speed should be compatible with edge-based processing needs. Additionally, CNN architectures are more conducive to pattern recognition (spatial fault prediction) for image data, whereas XGBoost is competitive for the structured numerical features used in predictive maintenance.

---

<sup>1</sup> Department of biomass, Al-Nahrain Research Center for Renewable Energy, Al-Nahrain University, Baghdad, Iraq

<sup>2</sup> Department of Wind Energy, Al-Nahrain Research Center for Renewable Energy, Al-Nahrain University, Baghdad, Iraq,

<sup>3</sup> Al-Mansour University College Computer Science and Information Systems Department

The presented framework is a practical option with a high sensitivity to photovoltaic system accuracy with a low computational cost. Its results are in favor of real-time operation and large-scale industrial deployment, which needs both detection performance and resource constraints, respectively.

**Keywords:** Photovoltaic systems, Fault diagnosis, Machine learning, Convolutional neural networks, Infrared thermography

## 1. Introduction

Photovoltaic (PV) energy systems are among the most widely deployed renewable technologies, driven by continuous improvements in conversion efficiency and sustained reductions in installation and operational costs [1], [2]. However, the rapid large-scale integration of PV installations has intensified the need for reliable operation to ensure stable and long-term power generation. PV modules are continuously exposed to harsh environmental conditions—such as temperature variations, humidity, and mechanical stress—which can trigger degradation mechanisms including micro-cracks, hotspots, delamination, and potential-induced degradation (PID) [1]. These defects may lead to considerable power losses, accelerated ageing, and a reduced operational lifetime of PV systems.

To ensure safe operation, minimise maintenance costs, and maintain performance stability, early and accurate fault diagnosis is essential. Conventional diagnostic techniques, including infrared (IR) thermography, electroluminescence (EL) imaging, and current–voltage (I–V) curve analysis, are widely employed for PV condition assessment [1], [2]. Nevertheless, their adoption in large-scale PV plants remains constrained by environmental sensitivity, reliance on expert interpretation, limited automation, and insufficient capability for real-time monitoring [2].

Recent advances in machine learning and deep learning have enabled data-driven diagnostic frameworks capable of processing large-scale, high-frequency datasets from imaging and electrical sensors [3], [4]. Accordingly, this study investigates the practical applicability of machine learning-based techniques for improving PV fault diagnosis, with a particular focus on real-time monitoring requirements and scalable deployment in industrial PV environments.

## 2. Related Work

### 2.1 Classical Fault Detection for PV Systems

Traditional fault detection techniques for photovoltaic (PV) systems rely on electrical parameters, infrared (IR) thermography, and electroluminescence (EL) imaging. Current–voltage (I–V) curve analysis techniques, including PV-based solutions, can effectively identify performance degradation; however, they provide limited spatial localisation of faults and are generally unsuitable for reliable real-time monitoring in large-scale PV installations [13], [14].

Infrared thermography is widely applied for diagnosing hotspots and overheating phenomena, yet its reliability is strongly influenced by environmental conditions, sensor calibration accuracy, and irradiance levels. Consequently, IR-based fault diagnosis often requires additional verification and expert interpretation, limiting its robustness in real-world operating environments [15]–[17].

Electroluminescence imaging enables high-resolution detection of micro-cracks and material-level defects through controlled current injection. However, since EL inspection is typically performed offline and requires specialised equipment, it lacks the practicality required for continuous and large-scale PV monitoring [19], [20]. These limitations collectively highlight the need for automated, scalable, and real-time fault detection methods suitable for modern photovoltaic systems.

## 2.2 PV Fault Detection Using Machine Learning-Based Methods

It is also of great interest that Machine learning (ML) techniques have gradually been used to automate PV fault detection and improve diagnostic performance [22], [23]. Deep learning algorithms especially Convolutional Neural Networks (CNNs) have shown promising results for analyzing EL and IR pictures through learning hierarchical features directly from dataset, being superior to classical handcrafted feature-based approaches [1], [2]. Classical ML methods such as Support Vector Machines (SVMs) and Extreme Gradient Boosting (XGBoost) have also been adapted for PV diagnostics such as structured electrical information analysis. Although it is computationally efficient, these methods generally do not perform as well on high-dimensional imaging data as CNN-based models [3]–[6]. Hybrid frameworks that combine CNN-based feature extraction with classical ML classifiers have, consequently, been proposed to maximize robustness [7], [8].

## 2.3 Limitations and Research Gap

Although recent breakthroughs have been made, these ML-based PV fault detection techniques all exhibit three main limitations. First, the use of one data modality decreases the robustness under environmental variation. Secondly, computational limitations frequently prevent the real-time and edge-level application [9], [10]. Third, generalizing the model across diverse PV technologies and environmental conditions is still an enormous challenge, which keeps research on adaptive/transfer-learning-based frameworks in the forefront [11]–[13].

## 2.4 Positioning of this Work

To overcome these limitations, we introduce a real-time UAV-based thermal fault detection framework explicitly accommodating environmental variations and deployment limitations. The approach combines deep learning and an optimized inference pipeline to achieve robustness and scalability for the scaling of its solution to large-scale PV monitoring in realistic operating scenarios

### 3. Problem Statement

As a primary renewable energy source, photovoltaic (PV) systems keep increasing in global adoption accompanied with new challenges for assuring their reliability and performance. Because PV modules are susceptible to many kinds of faults, such as microcracks, hotspots, delamination, and potential induced degradation, they can significantly decrease the energy output and operational lifetime. However, electrical parameter analysis, infrared (IR) thermography and electroluminescence (EL) imaging, etc. as the traditional fault detection methods, have many shortcomings. Typically, these methods do not scale well, are equipment specific, and often do not supply true real-time diagnostics, which restricts them from being employed in large scale PV farms. At the same time, with the rapid development of the research in machine learning (ML), it is possible to construct automated, accurate and scalable fault detection systems for PV modules. Unfortunately though, most ML based approaches introduced in previous studies center around a single algorithm, for instance, Convolutional Neural Networks (CNNs) or Support Vector Machines (SVMs), without fully having the ability to consider combined models with several learning techniques that result in better accuracy and reliability. In addition, the last frontier still includes real time monitoring capabilities and the ability to generalize among different PV environments. In this paper, we address these limitations by developing an advanced ML-based fault detection framework using deep learning (CNNs), traditional ML methods (SVMs, XGBoost) and hybrid methodologies to improve the fault classification accuracy, scalability and real time applicability. It aims to increase computational efficiency with high detection accuracy and allow PV systems to have more efficient performance at minimal maintenance costs. This study aims to bridge this gap as a means of making a better alternative to the existing fault detection methods for PV modules and promote the adoption of the sustainable solar energy as the broader solution to the environmental issues.

#### 3.1 Research Gap

Though a large number of works have attempted fault diagnosis in photovoltaic (PV) systems, several shortcomings still need to be taken care of in the existing works. Much of the existing literature studies a single diagnosis method like electroluminescence (EL) imaging, infrared thermography, or electrical current–voltage (I–V) measurements on their own. Although these methods have shown high performance in a controlled experimental scenario, they are generally not as robust and scalable in large scale PV system environments as in high power energy system settings. Second, studies often assess deep learning and classical machine learning models with various datasets or experimental setups, which severely reduces the validity of comparison and restricts the ability to determine the practical tradeoff between accuracy, computational cost, and deployment

flexibility. Moreover, insufficient attention has been paid to the applicability of fault diagnosis strategies for real-time or edge computing environments, where resource constraint means inference time and the usage of energy are the restrictions. Therefore, no unified research is available with systematic framework for the integration of multi modal PV diagnostic, easy and consistent comparison of different learning paradigms, as well as considering the real-time implementation requirements. Overcoming this divide is necessary to close the gap between high accuracy fault detection techniques and their realization at industrial-scale PV monitoring systems.

### 3.2 Research Objectives

To address the research gaps, this study aims to develop and evaluate a unified this paper proposes a comprehensive system fault diagnosis framework for photovoltaic systems.

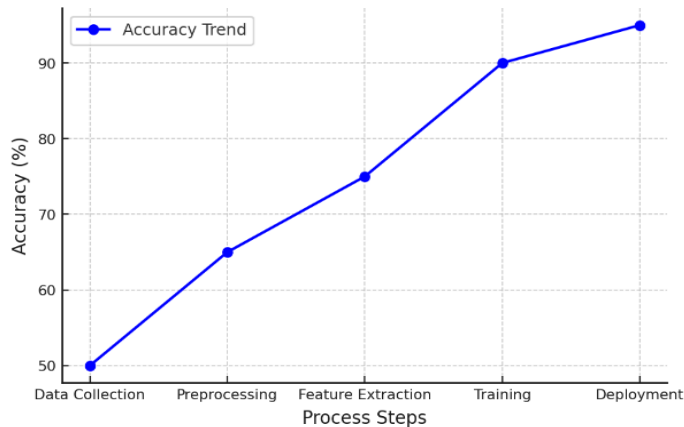
The objectives of this work are threefold:

- (i) to integrate electroluminescence imaging, infrared thermal measurements, and electrical I–V data within a single diagnostic framework;
- (ii) to provide a fair and systematic comparison between deep learning and classical machine learning models under identical experimental conditions; and
- (iii) to assess the suitability of the proposed approach for real-time and edge-based deployment, considering both detection accuracy and computational efficiency.

## 4. Methodology

In this section, the structured methodology for the development and evaluation of a machine learning-based photovoltaic (PV) fault detection system is presented. The proposed approach follows a systematic pipeline that includes data acquisition, data preprocessing, feature extraction, model training, evaluation, and real-time deployment, with the objective of maximising detection performance, scalability, and system reliability.

The fault detection framework is organised as an integrated data acquisition, processing, and real-time deployment pipeline, where each stage contributes to enhancing model accuracy and efficient fault diagnosis. Figure 1 illustrates the step-by-step workflow of the proposed system and demonstrates how preprocessing and feature extraction play a critical role in improving classification performance.



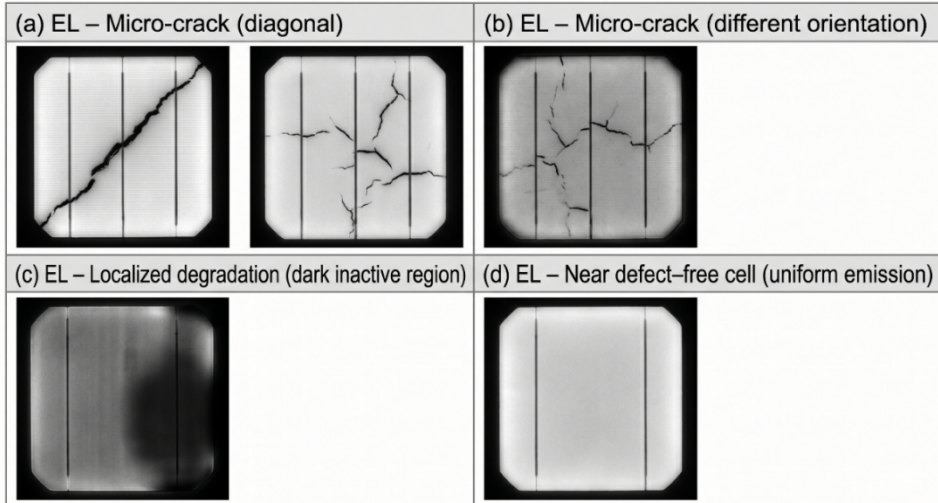
**Figure 1: Conceptual accuracy trend across the main stages of the proposed machine learning–based photovoltaic fault detection workflow**

In the final deployment stage, the trained model is implemented for real-time fault detection, enabling AI-based PV monitoring to serve as a practical and effective alternative to conventional diagnostic methods.

#### 4.1 Data Collection and Preprocessing

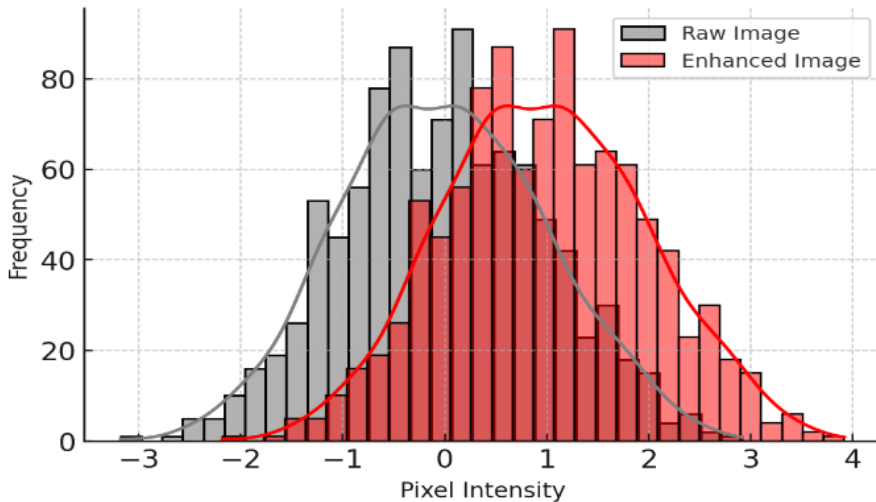
The data-driven fault detection depends fundamentally on good data quality and a good preprocessing pipeline. In this paper, high-resolution cameras were employed for image acquisition with electroluminescence (EL) and infrared (IR) images of photovoltaic (PV) cells (resolution  $<1024 \times 1024$  pixels). That was enough to identify the distinct defect features, such as micro-cracks, localised inactive areas, and structural damage. EL imaging was used to visualize electrically inactive zones and structural discontinuities at the cell level, and IR imaging to monitor thermal disturbances in the presence of non-uniform current flow and abnormal power output. Representative EL images showing different fault profiles such as diagonal and non-diagonal micro-cracks, local degradation regions, near defect-free cells and homogenous emission patterns are shown in Figure 2. Furthermore, in addition to imaging data for performance characterization the I–V characteristics of the system was measured on an environmental level in order to mimic realistic operating conditions for outdoor PV systems. Environment temperature tests were performed at  $-10\text{ }^{\circ}\text{C}$  up to  $60\text{ }^{\circ}\text{C}$ , with a solar irradiance of 200, 500, 800 and  $1000\text{ W/m}^2$ . This operating envelope allows us to study PV performance under common thermal and radiative stresses in field deployments. All EL and IR images were inspected and annotated manually per predefined fault categories, including micro-cracks, hotspots and local degradation. As depicted in Figure 2, the variety of observed defect manifestations provides a rationale for the development of machine learning

models that could learn complex spatial fault patterns and even adapt to heterogeneous PV fault patterns.



**Figure 2: Representative electroluminescence (EL) images of photovoltaic cells**

An image preprocessing pipeline, applied explicitly for both EL and IR images, has been shown to improve image quality and suppress noise on images, which could adversely affect model performance. First, undesired artefacts were preprocessed through  $3 \times 3$  Gaussian filtering. Adaptive histogram equalization was applied to refine contrast perception under different imaging regimes and clarify defect areas, such as micro-cracks and hotspots. The presence of structural defects, and robust feature extraction, was amplified by using Sobel and Canny edge detection filters. Furthermore, for overcoming overfitting and generalizing learning, the data augmentation technique was used (during training) such as rotation ( $\pm 30^\circ$ ), horizontal and vertical side flipping, brightness normalization and synthetic variation production among others. The combined impact of preprocessing on pixel-intensity distribution and feature performance is depicted in Figure 3.



**Figure 3: Preprocessing Techniques Applied to EL and IR Images**

This is evident in Figure 3, as preprocessing still effectively improves the contrast and clearness of image features. It is clear from the processed image that it exhibits a definite shift in pixel intensity distribution, hence have the visible defects. Such improvements allow EL and IR images to contain meaningful features from which machine learning models can classify fault (i.e. fault classification is more accurate).

To support image analysis, numerical transformation of the I-V curve data was performed to generate useful information. Cell degradation and power fluctuation pattern were analyzed using Wavelet Transform (WT) and Fourier Transform (FT). Several critical parameters as the mean voltage, the power deviation, the short circuit current and the standard deviation of irradiance response, were extracted by applying feature engineering techniques. Min max normalization technique was applied to standardize the data and to guarantee uniform feature distribution so that environmental fluctuation does not have a huge impact on the dataset. The dataset used in this study is publicly available and was originally published by Lazzaretti et al. in Sensors (2020). The data are used for non-commercial academic research in accordance with the Creative Commons BY-NC-SA 4.0 license.[24]

#### 4.2 Feature Extraction and Selection

Feature extraction and selection are very important steps to optimize the efficiency of machine learning models as only the most informative attributes can be used in classification and reducing computational overhead.

In order to analyze EL and IR images, deep learning models including Convolutional Neural Networks (CNN) were used for hierarchical feature

extraction. In the CNN architecture, five convolutional layers with ReLU activation were applied to capture features from low to high level such as basic edges and complex defect structure. To achieve a reduced computational complexity while preserving important defect information, max pooling layers with a  $2 \times 2$  kernel were implemented to operate on spatial dimensions. To bring in the activation values to a standardized value, the batch normalization layers were added that helped with the convergence speed of the model and with the stability of the model. To combat overfitting and achieve generalizability across different PV environments, dropout regularization of the 0.5 probability was used to deactivate random neurons during training.

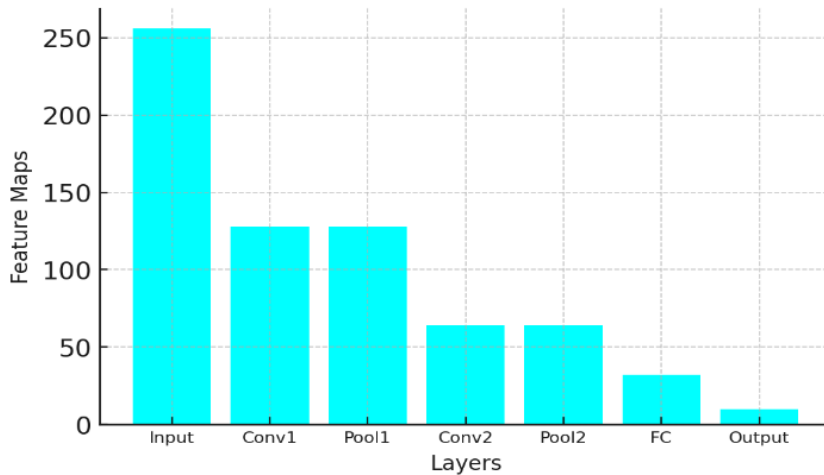
Multiple feature selection techniques were used in order to perform numerical research of I-V curve characteristics using model parameters. To address large dimensionality of the dataset and improve runtime efficiency, Principal Component Analysis (PCA) was employed to reduce dimensionality while maintaining 95% of the dataset's variance under the assumption that each component is orthogonal and uncorrelated to the others. We used Recursive Feature Elimination (RFE) to facilitate an iterative removal of less significant features, improving the dataset and highlighting the most important features for fault classification. To remove the redundant features and make sure the selected attributes contribute independently to the final model predictions, correlation analysis was performed.

### 4.3 Model Development and Training

To maximize the fault detection accuracy, the development of the hybrid machine learning framework was pursued by fusing the deep learning and the traditional machine learning.

To utilize various learning techniques with the better accuracy and robustness, a multi model fault detection pipeline was designed. Into this the CNNs based deep learning model was implemented to classify EL and IR images in that it is capable of automatic extraction of hierarchical image features. Multiple binary classification tasks of I-V curve features, namely, defective and non-defective modules, were realized using a Support Vector Machine (SVM). Finally, Extreme Gradient Boosting (XGBoost) was used for multi class classification of structured numerical data to improve predictive maintenance capabilities.

For the PV fault detection of CNN model, hierarchical feature extraction has been structured to perform in structured architecture. Each layer of this process helps to recognize the defects inside EL and IR images. Figure 4 presents the CNN model architecture with feature map sizes shown at stages in the CNN model



**Figure 4: CNN Model Architecture**

The model starts with an input layer and gangs of convolutional and pooling layers which progressively decrease spatial dimensions extracting complex features as seen in Figure 5. The output layer classifies the faults and Feature representations are further refined by fully connected layers (FC). Hierarchical feature learning on this architecture helps to effectively learn about defect in PV modules.

A supervised learning framework was then used to train the models given that 80% of the data was dedicated to the training and 20% to the validation. For the CNN model, the weight updates and the convergence speed were improved with Adam optimizer and a learning rate of 0.001 ( $\beta_1=0.9$ ,  $\beta_2=0.999$ ). The classification boundaries were optimized for CNN, XGBoost and for SVM the loss functions used were cross entropy loss for CNN and XGBoost and hinge loss for SVM.

For example, batch sizes of 32 images and for 512 samples have been used to ensure model efficiency for CNN and numerical models, respectively. To avoid over fitting, early stopping was applied to stop the training when validation loss did not improve for 10 consecutive epochs. To regularize the CNN layers, an L2 weight decay ( $\lambda=0.01$ ) was used. To achieve maximum performance, Grid Search and Random Search was applied to the extreme hyper parameter learning rates, dropout probabilities, and filter sizes to tune these hyper parameters for the best performance.

#### 4.4 Model Evaluation and Performance Metrics

In order to know that the PV fault detection performance in real world scenarios is reliable and accurate, evaluating the model performance is of most importance. Several metrics for KPI evaluation were analyzed:

- Percentage of correctly classified instances across all the fault categories, usually called 'Accuracy'.

- They measure it using Precision and Recall: Precision and Recall have to do with class-wise definitions of how effective a given detector is at defect detection, and how low can one make the number of false positives for a given coverage of potential defect sites.
- F1-Score: It is a balanced metric which takes into account precision and recall for a global performance assessment.
- ROC-AUC: Performance evaluation of models in ordering fault types.
- An insight is provided to misclassification rates and model errors by using the confusion matrix.
- Inferring Speed: Measured the real time classification capabilities of the deployed models.

#### **4.5 Real-Time Deployment and Testing**

Finally, to further verify the realism of the proposed marked based NIR inspection system, the trained models were deployed on edge computing devices to achieve efficient and low latency PV module diagnostic. To ensure real time fault detection, inference devices with less than 150ms per image were used to perform it, such as NVIDIA Jetson Nano and Raspberry Pi 4. We made use of it with automation fault detection alerting and decreased the on-site inspection with cloud integration on AWS IoT Core for remote the asset. Field test was made on a 5MW solar farm and the identified faults were detected in less than 2 second and per one full system scan. Deployment results also showed that the AI fault detection system can be used in a scalable and cost-effective manner for observing PV monitoring.

#### **4.6 Results and Discussion**

The subsequent section provides an extensive evaluation and model assessment of the proposed machine learning based photovoltaic (PV) fault detection approach. Model performance is assessed in terms of classification accuracy, precision, recall, F1 score, false positive rate (FPR), false negative rate (FNR), and inference latency. Also the effect of feature selection strategies, real time deployment feasibility and scalability in practical operating conditions is analysed.

#### **4.7 Model Performance Comparison**

In Table 2 summarization are shown the performance of the machine learning models evaluated. The results demonstrate that CNN-based model performance clearly outperforms SVM and XGBoost with respect to the evaluation indices. Both its classification accuracy, reaching up to 95.2%, the CNN represents the highest and the very best precision with the highest recall, which signifies the balanced detection of both defective and non-defect PV samples. The CNN exhibits the lowest false-negative rate (2.1%), which is especially salient for safety-based PV monitoring, where ignored faults could result in energy losses or system damage.

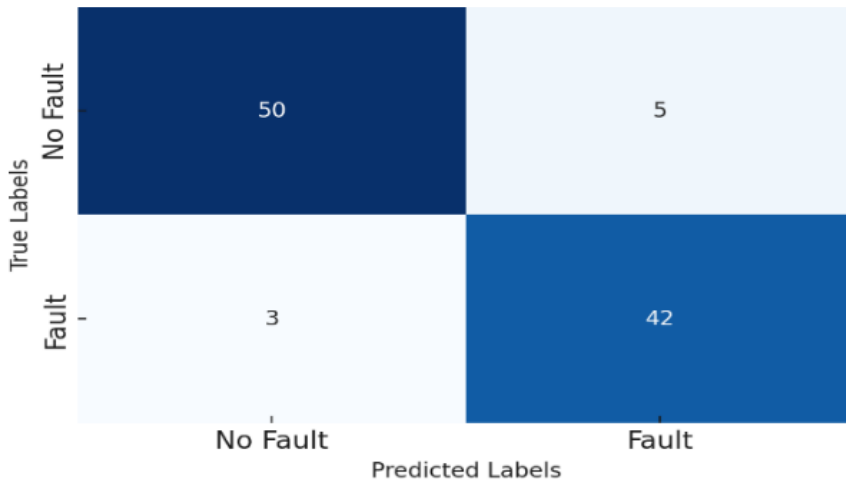
The quality (92.3% accuracy) of structured numerical features which depend on the I–V characteristics are also good performance of XGBoost, proving its effectiveness not only as well-accepted for degradation trending and for predictive maintenance algorithms. On the contrary, it has a relatively high false-negative rate (4.6%) and a low capability of high-dimensional image data modeling, which limit its effectiveness for image-based faults detection. The SVM model achieves the lowest level of performance (88.5% accuracy), with higher rate of misclassification, indicating more dependence on handcrafted features, and less scalability of the system to deal with complex PV defect patterns. From a real-time perspective, all models perform within near real-time operational requirements. Due to the greater latency (140 ms versus SVM, 90 ms), however, its accuracy and reliability are much better than SVM, hence the higher latency is required for the CNN. In general, the results indicate there is the best trade-off between detection performance and deployability for CNN, for which XGBoost is much better as a replacement for numerical-driven maintenance analysis.

**Table 1. Performance comparison of machine learning models for photovoltaic fault detection in terms of accuracy, classification metrics, and inference time**

| Model   | Accuracy (%)          | Precision (%) | Recall (%) | F1-score (%) | False Positive Rate (FPR) | False Negative Rate (FNR) | Inference Time (ms) |
|---------|-----------------------|---------------|------------|--------------|---------------------------|---------------------------|---------------------|
| CNN     | 95.2<br>( $\pm 1.3$ ) | 96.1          | 94.5       | 95.3         | 1.8%                      | 2.1%                      | 140                 |
| XGBoost | 92.3<br>( $\pm 1.7$ ) | 93.0          | 91.2       | 92.1         | 3.5%                      | 4.6%                      | 110                 |
| SVM     | 88.5<br>( $\pm 2.1$ ) | 89.1          | 86.7       | 87.9         | 5.3%                      | 6.8%                      | 90                  |

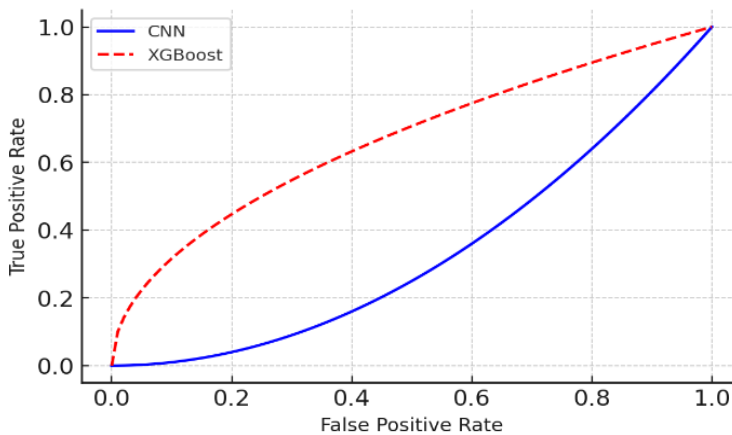
#### 4.8 Confusion Matrix and ROC Analysis

The confusion matrix of the CNN-based model is shown in Figure 5, which indicates a high proportion of correctly classified fault samples and a low false-alarm rate.



**Figure 5: Confusion Matrix for Model Performance Evaluation**

Only a small number of fault samples are misclassified as normal, indicating strong reliability under realistic conditions. Receiver Operating Characteristic (ROC) analysis was also used to test how well the models discriminatively perform under the different decision thresholds applied. Figure 6 shows that the CNN has a high AUC value that indicates good separation between the faulty and non-faulty PV samples.



**Figure 6: Receiver Operating Characteristic (ROC) curve comparison between CNN and XGBoost models for PV fault classification.**

Although XGBoost is competitive for ROC for structured numerical features, the CNN exhibits better practical discrimination for image-based diagnostics especially with spatially complex faults at EL images and IR images. These findings suggest that ROC performance does not fully represent real-world effectiveness; CNNs

provide a higher operational reliability due to greater accuracy and lower FNR, requirements of a field-deployed PV fault detection system.

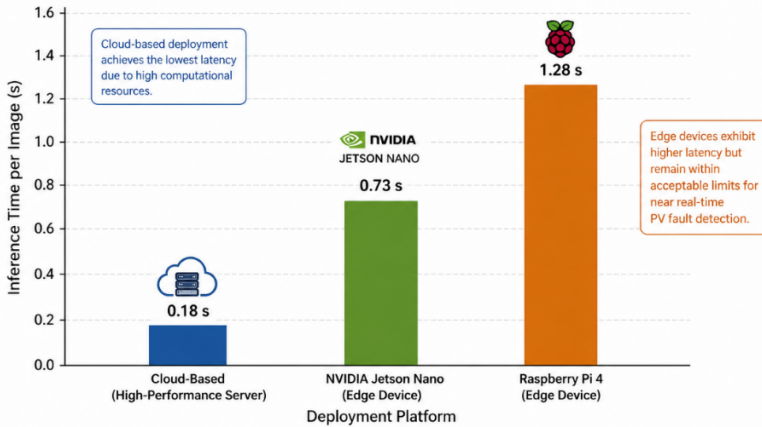
#### **4.9 Impact of Feature Selection on Performance**

Feature selection techniques were employed to evaluate their impact on detection accuracy and computational efficiency. Principal Component Analysis (PCA) preserved approximately 95% of the variance in the original feature space, resulting in a 23% reduction in training time and a modest improvement in classification accuracy from 94.8% to 95.2%. This indicates that PCA is an effective approach for eliminating feature redundancy while retaining the most discriminative characteristics.

In addition, Recursive Feature Elimination (RFE) improved model generalisation by removing less informative features, leading to a 19% reduction in training time and an increase in classification accuracy from 91.1% to 92.3%. These results collectively demonstrate that feature selection plays a critical role in achieving an optimal trade-off between predictive accuracy and computational cost, which is particularly important for real-time and edge-based photovoltaic monitoring applications.

#### **4.10 Scalability and Real-Time Deployment**

To find real time feasibility, real-time feasibility of the proposed framework was validated by executing on NVIDIA Jetson Nano and Raspberry Pi 4. The inference latency for the CNN model was less than 150 ms per image, fulfilling the near real-time operational needs. At a 5 MW solar farm, field evaluation demonstrated scalability, with diagnostic latency for full system below the 2 s per scan range. Integrating with AWS IoT Core allowed for hybrid edge–cloud monitoring with real-time fault reporting, centralised supervision with reduced dependency on manual inspections. In Figure 7, it is illustrated that cloud-based deployment works very quickly for inference but requires strong connectivity, while edge-based deployment is more robust and autonomous for rural PV sites.



**Figure 7 Comparison of inference latency for CNN-based photovoltaic fault detection across different deployment platforms**

#### 4.11 Future Developments in AI for PV Monitoring

Feature selection was investigated to see the performance and computational efficiency effects of the model. Similar approaches to the detection of photovoltaic fault; Principal Component Analysis (PCA) and Recursive Feature Elimination (RFE) were carried out to optimize the feature space without compromising the discriminative information required for failure detection. According to Table 3, PCA takes up around 95% of the initial variance in the data, reducing feature dimensionality significantly. Also achieved with this reduction was the 23% decrease in training time and the improvement in classification accuracy from 94.8% to 95.2%. Results indicate a large portion of features are removed without having to drop, as a result, that PCA is highly efficient at removing from the dataset unnecessary correlations which allows the learning model to learn toward the accuracy of prediction in a reasonable interval of time. RFE also attained better model accuracy by de-constructing less informative features step by step. Using RFE the training time was optimized by 19% and the classification accuracy was improved from 91.1% to 92.3%, see Table 2. RFE enhances the generalization in the model and decreases overfitting, particularly in high dimensional numerical feature series by focusing on sensitive features. We show feature selection is essential for achieving the appropriate accuracy-to-computation cost trade-off. PCA and RFE combined render cost-effective and scalable PV fault detection, which is essential for near real time monitoring and edge-based deployment in large-scale photovoltaic systems.

**Table2.Impact of Feature Selection on Model Performance.**

| Feature Selection Method | Training Time Reduction (%) | Accuracy Before (%) | Accuracy After (%) |
|--------------------------|-----------------------------|---------------------|--------------------|
| PCA (Top 20 Components)  | 23%                         | 94.8                | 95.2               |
| RFE (Feature Pruning)    | 19%                         | 91.1                | 92.3               |

## 5 Conclusion and Future Work

The rapid implementation of photovoltaic (PV) energy systems needs fault detection methods that can be accurately developed in real-time and scale fast. This study shows that machine learning-based fault detection is an effective and simple substitute for classical fault diagnosis for large PV systems. Through analyzing of the experiments, it is found that CNN received the greatest performance on fault classification as it reached the best accuracy value of 95.2% and low false negative rate of 2.1%, through the extraction of the discriminative spatial features from EL and IR visualizations. While on structured (numerical) data and predictive maintenance, XGBoost was claimed to have been an effective choice for the latter, Support Vector Machines (SVMs), while computationally efficient, were noted to generate higher false-positive rates than the former and thus limited their practical use in safety critical scenarios. Together PCA and RFE for feature selection yielded low computational costs and detection quality. Furthermore, inference latency was confirmed under 150 ms for the deployment, on edge environments with devices like NVIDIA Jetson Nano and Raspberry Pi architectures, and the scalable and reliable framework was tested in real field settings: total system diagnostics latency is below 2 s on a 5 MW solar farm. Yet there will always be some limitations. The dataset used here can be limited in geographical context and much more data will need to be collected across different climates, PV technologies, and aged conditions to enhance the ability of the model to generalize. Advanced network lightweight CNN models such as MobileNet and quantized neural network architectures will be further optimized to maximize the performance of edge-devices. In addition, to make them practical in industry at large, we will need to enhance the interpretability of the models with explainable AI (XAI) approaches like Grad-CAM, SHAP and LIME. Additional research on federated learning for privacy-aware distributed training over multiple PV sites, transformer-based vision models, and time series learning for enhancing fault representation and preventive maintenance using various sites will be further investigated. These advances will facilitate autonomous, resilient, and sustainable AI-based PV monitoring solutions.

## 6. References

- [1] L. Zhang, X. Wu, Z. Liu, P. Yu, and M. Yang, "ESD-YOLOv8: An efficient solar cell fault detection model based on YOLOv8," *IEEE Access*, vol. 12, pp. 138801–138815, 2024, doi: 10.1109/ACCESS.2024.3466209.
- [2] I. U. Khalil et al., "Comparative analysis of photovoltaic faults and performance evaluation of its detection techniques," *IEEE Access*, vol. 8, pp. 26676–27000, 2020, doi: 10.1109/ACCESS.2020.2970531.
- [3] X. Zhang, T. Hou, Y. Hao, H. Shangguan, A. Wang, and S. Peng, "Surface defect detection of solar cells based on multiscale region proposal fusion network," *IEEE Access*, vol. 9, pp. 62093–62101, 2021, doi: 10.1109/ACCESS.2021.3074219.
- [4] S. Lu, K. Wu, and J. Chen, "Solar cell surface defect detection based on optimized YOLOv5," *IEEE Access*, vol. 11, pp. 71026–71036, 2023, doi: 10.1109/ACCESS.2023.3294344.
- [5] M. Lavador-Osorio et al., "An enhanced frequency analysis and machine learning-based approach for open circuit failures in PV systems," *IEEE Access*, vol. 12, pp. 96342–96357, 2024, doi: 10.1109/ACCESS.2024.3425486.
- [6] M. Abdelsattar et al., "Automated defect detection in solar cell images using deep learning algorithms," *IEEE Access*, vol. 13, pp. 4136–4157, 2025, doi: 10.1109/ACCESS.2024.3525183.
- [7] M. Zhang and L. Yin, "Solar cell surface defect detection based on improved YOLOv5," *IEEE Access*, vol. 10, pp. 80804–80815, 2022, doi: 10.1109/ACCESS.2022.3195901.
- [8] B. Xia et al., "Surface defect recognition of solar panel based on percolation-based image processing and Serre standard model," *IEEE Access*, vol. 11, pp. 55126–55138, 2023, doi: 10.1109/ACCESS.2023.3281653.
- [9] S. Datta et al., "A comprehensive review of the application of machine learning in fabrication and implementation of photovoltaic systems," *IEEE Access*, vol. 11, pp. 77750–77778, 2023, doi: 10.1109/ACCESS.2023.3298542.
- [10] P.-Y. Sevilla-Camacho et al., "A novel fault detection and location method for PV arrays based on frequency analysis," *IEEE Access*, vol. 7, pp. 72050–72061, 2019, doi: 10.1109/ACCESS.2019.2920053.
- [11] S. P. Pathak and S. A. Patil, "Evaluation of effect of pre-processing techniques in solar panel fault detection," *IEEE Access*, vol. 11, pp. 72848–72860, 2023, doi: 10.1109/ACCESS.2023.3293756.

- [12] J. Huang et al., "A novel MoCo-based self-supervised learning framework for solar panel defect detection," *IEEE Access*, vol. 13, pp. 22977–22988, 2025, doi: 10.1109/ACCESS.2025.3529701.
- [13] H. Kang et al., "Photovoltaic cell defect detection based on weakly supervised learning with module-level annotations," *IEEE Access*, vol. 12, pp. 5575–5583, 2024, doi: 10.1109/ACCESS.2024.3349975.
- [14] H. Tan et al., "RAFBSD: An efficient detector for accurate identification of defects in photovoltaic cells," *IEEE Access*, vol. 12, pp. 61512–61528, 2024, doi: 10.1109/ACCESS.2024.3393934.
- [15] G. G. Kim et al., "Fault detection for photovoltaic systems using multivariate analysis with electrical and environmental variables," *IEEE J. Photovolt.*, vol. 11, no. 1, pp. 202–212, 2021, doi: 10.1109/JPHOTOV.2020.3032974.
- [16] F. M. A. Mazen et al., "Deep learning for automatic defect detection in PV modules using electroluminescence images," *IEEE Access*, vol. 11, pp. 57783–57795, 2023, doi: 10.1109/ACCESS.2023.3284043.
- [17] P. A. A. Pramana and R. Dalimi, "Photovoltaic thermal fault monitoring using the catadioptric device," *IEEE Access*, vol. 11, pp. 75546–75554, 2023, doi: 10.1109/ACCESS.2022.3203814.
- [18] Z. Wu et al., "A review for solar panel fire accident prevention in large-scale PV applications," *IEEE Access*, vol. 8, pp. 132466–132480, 2020, doi: 10.1109/ACCESS.2020.3010212.
- [19] D. P. Winston et al., "Solar PV micro-crack and hotspot detection techniques using neural networks and SVM," *IEEE Access*, vol. 9, pp. 127259–127269, 2021, doi: 10.1109/ACCESS.2021.3111904.
- [20] A. Procházka et al., "Advanced signal processing techniques for monitoring east/west oriented solar photovoltaic systems," *IEEE Access*, vol. 12, pp. 165042–165049, 2024, doi: 10.1109/ACCESS.2024.3492017.
- [21] Z. Shen et al., "Defect detection in c-Si photovoltaic modules via transient thermography and deconvolution optimization," *Chinese J. Electr. Eng.*, vol. 10, no. 1, pp. 3–11, 2024, doi: 10.23919/CJEE.2023.000043.
- [22] H. P.-C. Hwang et al., "Detection of malfunctioning photovoltaic modules based on machine learning algorithms," *IEEE Access*, vol. 9, pp. 37210–37219, 2021, doi: 10.1109/ACCESS.2021.3063461.

- [23] M. Hussain et al., “A gradient guided architecture coupled with filter-fused representations for micro-crack detection in photovoltaic cells,” IEEE Access, vol. 10, pp. 58950–58964, 2022, doi: 10.1109/ACCESS.2022.3178588.
- [24] A. E. Lazzaretti et al., “Monitoring system for online fault detection and classification in photovoltaic plants,” Sensors, vol. 20, no. 17, pp. 1–23, 2020.

## تقييم خوارزميات التعلّم الآلي لاكتشاف الأعطال في الخلايا الشمسية: الأداء والقيود

م . رافت عبد الكاظم عبيدة<sup>2</sup>[raafatoubida@nahrainuniv.edu.iq](mailto:raafatoubida@nahrainuniv.edu.iq)م . م . زاهر فاضل رحم<sup>1</sup>[Zahir.raham@gmail.com](mailto:Zahir.raham@gmail.com)أ.د. عبد المنعم صالح رحمة<sup>3</sup>[monem.rahma@muc.edu.iq](mailto:monem.rahma@muc.edu.iq)

**المستخلص:** شهدت أنظمة الطاقة الكهروضوئية (PV) توسعًا سريعًا في استخدامها، مما أدى إلى زيادة الحاجة إلى أساليب كشف الأعطال التي تتسم بالدقة، وانخفاض التكلفة، والعمل في الزمن الحقيقي. وتُعد تقنيات التصوير الحراري بالأشعة تحت الحمراء، والتصوير بالإضاءة الكهربائية (Electroluminescence – EL)، وتحليل منحنيات التيار-الجهد (I–V) من أكثر الأساليب التشخيصية شيوعًا في أنظمة الخلايا الكهروضوئية، إلا أنه لا يزال من الصعب توسيع نطاقها وأتمنتها على مستوى الأنظمة الواسعة. في هذا البحث، نقتراح منهجية تعلّم آلي متعددة الوسائط لكشف أعطال أنظمة PV من خلال دمج صور EL، وقياسات الأشعة تحت الحمراء الحرارية، والخصائص الكهربائية لمنحنيات I–V. كما أجرينا مقارنات منهجية بين نماذج التعلّم العميق ونماذج التعلّم الآلي التقليدية، وهي: الشبكات العصبية الالتفافية (CNN)، وآلات المتجهات الداعمة (SVM)، وخوارزمية التعزيز التدريجي المتطرف (XGBoost)، باستخدام البنى البيانية نفسها. تم اعتماد معالجات متقدمة للصور واختيار فعال للخصائص من خلال دمج تحليل المكونات الرئيسية (PCA) وخوارزمية الحذف التراجعي للخصائص (RFE) لتقليل العبء الحاسوبي. وبالاعتماد على تقييم تجريبي لمجموعة بيانات حقيقية مكونة من 2,624 صورة EL، حقق نموذج CNN أعلى دقة تصنيف بلغت 95.2% وزمن استجابة أقل من 150 ميلي ثانية على أجهزة الحوسبة الطرفية مثل NVIDIA Jetson Nano. وتُظهر النتائج أن نموذج CNN أكثر كفاءة في كشف أعطال PV المعتمدة على الصور، في حين يُظهر XGBoost أداءً قوياً مع البيانات العددية المنظمة في تطبيقات الصيانة التنبؤية. ومن خلال التركيز على دقة الكشف والكفاءة وإمكانية التشغيل في الزمن الحقيقي، يسهم هذا العمل في سد الفجوة بين كشف أعطال أنظمة PV وأنظمة المراقبة الشمسية على المستوى الصناعي.

**الكلمات المفتاحية:** لطاقة الكهروضوئية (PV)، التصوير الحراري بالأشعة تحت الحمراء (IR)، التعلّم الآلي، الشبكات العصبية الالتفافية (CNN)، آلات المتجهات الداعمة (SVM)، التعزيز التدريجي المتطرف (XGBoost)، إدارة الطاقة الشمسية، التصوير بالإضاءة الكهربائية (EL)، كشف الأعطال، تحسين الصور، تقليل الضوضاء.

<sup>1</sup> جامعة النهريين مركز بحوث النهريين للطاقة المتجددة /قسم الكتلة الحيوية

<sup>2</sup> جامعة النهريين مركز بحوث النهريين للطاقة المتجددة/قسم طاقة الرياح

<sup>3</sup> كلية المنصور الجامعة