# *Machine Learning–Based Performance Prediction and Health Assessment for Optoelectronic Devices Using Optical–Electrical Feature Fusion*

**Shaoyi Sun[1,a,#], Chunyu Ma[1,b,#]**

[1]*School of Electronic and Optical Engineering, Nanjing University of Science and Technology, Nanjing, China*
[#]*These authors contributed equally to this work*
[a]*15533188817@163.com,* [b]*17351387672@163.com*

*Keywords:* Optoelectronic devices; machine learning; optical–electrical feature fusion; performance prediction; health assessment; multi-task learning

*Abstract:* Optoelectronic devices such as light-emitting diodes, photodetectors, and small photovoltaic modules are widely used in lighting, sensing, and energy conversion. However, their performance and health status are strongly influenced by manufacturing variability, operating conditions, and gradual degradation, which makes traditional threshold-based evaluation methods inaccurate and labor-intensive. In this work, we propose a machine learning–based framework for performance prediction and health assessment of optoelectronic devices using fused optical–electrical features. First, optical spectra, luminous flux, and chromaticity coordinates are combined with electrical characteristics such as I–V curves, input power, and operating temperature to construct a comprehensive feature set. Principal component analysis and normalization are then applied to reduce redundancy and stabilize the input space. On top of these features, we systematically compare several regression and classification models, including Random Forest, gradient boosting–based methods, and deep neural networks. A multi-task learning strategy is further introduced to jointly predict key performance indicators (e.g., efficiency, output power) and discrete health states (healthy, mildly degraded, severely degraded), using a hybrid loss that balances regression accuracy and classification robustness. Experimental results on a mixed dataset of simulated and measured devices show that the proposed framework achieves higher coefficient of determination and lower error metrics than conventional single-feature or single-model baselines, while providing interpretable feature importance that is consistent with physical intuition. The study demonstrates the feasibility of integrating optical–electrical measurements with machine learning for intelligent monitoring and remaining-life assessment of optoelectronic devices.

## 1. Introduction

Optoelectronic devices such as light-emitting diodes (LEDs), organic LEDs (OLEDs), dye-sensitized solar cells (DSSCs), and photovoltaic (PV) modules underpin modern lighting, display,

and energy-conversion technologies, yet their performance and lifetime remain highly sensitive to material quality, device structure, and operating conditions. Traditional design and qualification workflows rely heavily on physics-based simulation, accelerated testing, and threshold-based pass/fail criteria, which are costly and often struggle to capture complex multi-parameter interactions. Recent studies have demonstrated that machine learning (ML) can effectively learn nonlinear mappings from device descriptors or process parameters to key performance metrics—for example, predicting DSSC power-conversion efficiency from electrical parameters and fabrication variables, or using active learning to optimise GaN-based LED structures and mitigate efficiency droop. [1,2] At the same time, ML and deep learning models have been integrated with compact device models and digital-twin concepts to enable prognostics and health management (PHM) of LEDs and OLEDs, providing data-driven estimates of remaining useful life directly from optical and electrical degradation signatures. [3,4] These developments highlight the potential of ML to move beyond simple curve-fitting and act as a core engine for intelligent optoelectronic device analytics.

Parallel advances have emerged on the system side, where ML-based condition monitoring and PHM frameworks are increasingly applied to PV arrays and other power-conversion assets under realistic field conditions. Berghout et al. reviewed ML-based condition monitoring for PV systems, showing how supervised algorithms can detect faults and performance losses from electrical operating data under varying irradiance and temperature profiles. [5] Jobayer et al. further provided a systematic review of ML methods for predicting PV system parameters, emphasising the importance of robust feature engineering, model selection, and uncertainty handling when dealing with heterogeneous measurement sources. [6] Despite this progress, most existing works still focus on either device-level optimisation or system-level monitoring in isolation, and often use predominantly electrical or predominantly optical inputs. There is comparatively little work that treats fused optical–electrical measurements as a unified feature space for both performance prediction and health assessment at the device level, within a single multi-task learning framework. This motivates the present study, which aims to develop and evaluate a machine learning–based framework that jointly exploits optical and electrical features for predicting key performance indicators and quantifying health states of optoelectronic devices.

## 2. Related Work

In recent years, machine learning has become a central tool for modelling structure–property relationships in optoelectronic materials and devices, especially organic and perovskite solar cells. Mahmood and Wang systematically reviewed how supervised learning and Bayesian optimisation accelerate donor–acceptor screening and morphology optimisation for high-efficiency organic solar cells, highlighting both data scarcity and generalisation issues in current models [7]. Meftahi et al. showed that kernel methods and ensemble learning can predict key device metrics (open-circuit voltage, short-circuit current and power conversion efficiency) from molecular descriptors and processing conditions with accuracy sufficient to guide experimental screening [8]. Building on these foundations, Ahmed and co-workers summarised a broad landscape of OSC applications, from data-driven molecular design to device lifetime prediction, and emphasised the need for interpretable features that reflect underlying photophysical mechanisms rather than black-box correlations [9]. More recent work has therefore embedded explainability directly into the modelling pipeline: Siddiqui et al. combined SHAP analysis with classification models to identify donor–acceptor pairs that consistently yield high-efficiency devices [10], while Lee et al. designed an interpretable model for ternary OSCs that reveals how the third component modifies charge transport and recombination pathways [11]. Rodrigues et al. adopted a data-centric perspective,

using careful curation and augmentation to stabilise performance prediction across diverse organic semiconductors [12]. Beyond photovoltaic materials, similar ideas have been applied to light-emitting devices: Yuan et al. used deep neural networks to learn the spectral power distribution of LEDs under multiple degradation mechanisms, effectively turning time-resolved optical measurements into predictors of ageing behaviour [13], and Błaszczak and Gryko analysed multi-emitter LED retrofits with statistics-aware metrics (such as the reliability of $R^2$) to quantify fit quality and guide model selection for spectral reconstruction [14]. Together, these studies demonstrate that machine learning can capture complex couplings between optical spectra, device structure and performance, but they also reveal limitations such as weak physical constraints, limited use of multi-modal data, and insufficient attention to uncertainty quantification.

At the system level, research on photovoltaic (PV) modules and arrays has driven another major branch of optoelectronic machine learning focused on defect diagnosis and prognostics. For image-based module inspection, Masita et al. reviewed deep learning methods for electroluminescence, infrared and RGB imagery, noting that convolutional networks now routinely detect cracks, hotspots and soiling but often remain tuned to a single dataset or fault family [15]. Et-taleby et al. surveyed machine-learning algorithms for PV fault detection and compared electrical-signal-based and image-based schemes, concluding that artificial neural networks and CNNs dominate the literature but still struggle with generalisation across plants and operating conditions [16]. Islam et al. extended this perspective by analysing artificial-intelligence-driven PV fault identification and diagnosis, with particular emphasis on hybrid electrical/thermal monitoring and the role of data quality and labelling [17]. At the algorithmic level, Sabati et al. integrated multiple pre-trained CNNs with a bio-inspired Bitterling Fish Optimisation algorithm to improve feature selection and classification of PV panel defects from RGB images [18], while Nassreddine et al. designed ensemble tree-based models that distinguish normal operation from seven fault types under different operating modes, achieving near-perfect accuracy after hyperparameter tuning [19]. Hybrid and optimised deep-learning architectures have also been proposed for large-scale PV plants: Bougoffa et al. combined CNNs with recurrent and dense layers to diagnose multiple faults under varying irradiance and temperature [20], and Khandeparkar et al. systematically compared supervised classifiers (decision trees, random forests, SVMs, XGBoost) for electrical fault detection in grid-connected PV systems, showing that gradient-boosted trees often deliver the best trade-off between accuracy and robustness [21]. Teta et al. further demonstrated that lightweight CNNs tuned with a meta-heuristic Energy Valley Optimizer can be deployed on edge devices for real-time PV fault diagnosis [22]. Overall, the current state of the art confirms the effectiveness of machine learning for optical–electrical condition monitoring, but most studies still treat optical and electrical channels separately, focus on static datasets rather than continuous degradation trajectories, and rarely enforce physics-based consistency—leaving room for integrated, optics-aware frameworks that fuse multi-modal signals and embed device physics into the learning process.

## 3. Methods

### 3.1. Overall Framework

The proposed framework treats performance prediction and health assessment of optoelectronic devices as a supervised learning problem built on fused optical–electrical features. First, raw measurements from each device are acquired, including optical spectra, luminous flux, chromaticity coordinates, and electrical characteristics such as I–V curves, input power and operating temperature. These heterogeneous signals are then preprocessed via cleaning, normalization and dimensionality reduction to form a compact numerical feature vector for each device. Next, a set of supervised models is trained to map these features to target quantities: continuous performance

indicators (e.g., efficiency, output power) and discrete health states (healthy, mildly degraded, severely degraded). The models include both classical ensemble learners and deep neural networks, organized in a multi-task architecture that jointly optimizes regression and classification. Finally, trained models are evaluated using a held-out test set and cross-validation, and their behaviour is interpreted through feature importance scores and partial dependence analysis, allowing us to link data-driven predictions back to optical and electrical device physics. The implementation follows a modular pipeline design similar to that popularized in general-purpose machine learning toolkits such as scikit-learn, where preprocessing and learning stages are chained and fitted jointly to avoid data leakage between training and testing. [23]

Formally, consider a dataset of $N$ devices, each represented by a fused feature vector $\mathbf{x}_i \in \mathbb{R}^d$, a set of continuous performance labels $\mathbf{y}_i^{(\text{perf})} \in \mathbb{R}^K$ (e.g., efficiency, output power) and a health-state label $y_i^{(\text{health})} \in \{1, ..., C\}$. The goal is to learn a parametric predictor $f_\theta$ such that

$$(\hat{\mathbf{y}}_i^{(\text{perf})}, \hat{\mathbf{p}}_i^{(\text{health})}) = f_\theta(\mathbf{x}_i), \tag{1}$$

where $\hat{\mathbf{p}}_i^{(\text{health})}$ denotes the estimated class probabilities over health states. The parameters $\theta$ are optimized by minimizing a joint loss function that combines regression and classification objectives, as detailed below.

## 3.2. Feature Construction and Preprocessing

The raw measurements for each device can be partitioned into optical and electrical components. Let $\mathbf{s}_i \in \mathbb{R}^{d_s}$ denote the sampled emission spectrum (or other wavelength-resolved optical signal) of device $i$, and let $\mathbf{e}_i \in \mathbb{R}^{d_e}$ collect scalar optical and electrical features such as luminous flux $\Phi$, chromaticity coordinates $(x, y)$, characteristic points on the I–V curve, input power $P_{\text{in}}$ and operating temperature $T$. To stabilize the learning problem and reduce redundancy, we first standardize each scalar feature across the dataset as

$$\tilde{e}_{i,j} = \frac{e_{i,j} - \mu_j}{\sigma_j}, j = 1, ..., d_e, \tag{2}$$

where $\mu_j$ and $\sigma_j$ are the empirical mean and standard deviation of feature $j$ computed on the training set.

The high-dimensional spectral vector $\mathbf{s}_i$ is compressed using principal component analysis (PCA). We compute the empirical mean spectrum $\bar{\mathbf{s}}$ and the covariance matrix

$$\Sigma = \frac{1}{N-1} \sum_{i=1}^{N} (\mathbf{s}_i - \bar{\mathbf{s}})(\mathbf{s}_i - \bar{\mathbf{s}})^{\mathsf{T}}, \tag{3}$$

Then obtain its eigen-decomposition $\Sigma = W \Lambda W^{\mathsf{T}}$, where columns of $W$ are orthonormal eigenvectors. The spectrum $\mathbf{s}_i$ is projected onto the first $k$ principal components as

$$\mathbf{z}_i = W_k (\mathbf{s}_i - \bar{\mathbf{s}}) \in \mathbb{R}^k, \tag{4}$$

where $W_k$ contains the eigenvectors corresponding to the largest $k$ eigenvalues, chosen to explain at least 95% of the spectral variance. This step is motivated by recent work on degradation tracking and luminescence mapping of optoelectronic devices, where compressed optical representations retain the dominant physical information while dramatically reducing dimensionality and noise. [25]

Finally, the optical and electrical descriptors are concatenated to form the fused feature vector

$x_i = [z_i; \tilde{e}_i] \in \mathbb{R}^{k+d_e}$, which serves as the input to all downstream models.

Where appropriate, interaction features such as ratios $\Phi/P_{\text{in}}$ (luminous efficacy) or temperature coefficients $\Delta\Phi/\Delta T$ are added to encode simple, physically motivated relationships that can ease the burden on the learning algorithms.

## 3.3. Ensemble Learning for Performance Prediction

As a first family of models, we employ tree-based ensemble methods for predicting the continuous performance targets $\mathbf{y}_i^{(\text{perf})}$ from the fused features $\mathbf{x}_i$. Random forest (RF) regression constructs an ensemble of $M$ decision trees $\{T_m\}_{m=1}^{M}$, each trained on a bootstrap sample of the training data and a randomly selected subset of features at each split. The RF predictor is given by

$$\hat{\mathbf{y}}_i^{(\text{perf})} = \frac{1}{M}\sum_{m=1}^{M} T_m(\mathbf{x}_i),$$

(5)

where each $T_m$ outputs a vector of performance indicators. By aggregating many decorrelated trees, RF reduces variance and provides an inherent estimate of feature importance, which we later use to interpret the influence of optical versus electrical descriptors.

In addition to RF, we adopt gradient boosting regression, which builds the predictor as a stage-wise sum of weak learners (typically shallow trees). Starting from an initial guess $F_0(\mathbf{x})$, the model at iteration $t$ is updated as

$$F_t(\mathbf{x}) = F_{t-1}(\mathbf{x}) + \eta h_t(\mathbf{x}),$$

(6)

where $\eta \in (0,1]$ is the learning rate and $h_t$ is fit to the negative gradient of the loss function with respect to the current predictions. When applied to the mean-squared error loss, this procedure can be viewed as functional gradient descent in function space, leading to strong approximators that capture complex nonlinear interactions in the data. The success of such ensemble models in predicting electrical characteristics of PV modules from I–V curves and environmental variables motivates their use here as competitive baselines on our fused optical–electrical feature set. [26]

## 3.4. Deep Neural Networks for Joint Performance and Health Prediction

To more fully exploit subtle correlations between optical spectra, electrical behaviour and long-term degradation, we design a multi-task deep neural network (DNN) that shares a common feature extractor across tasks and branches into separate heads for regression and classification. The shared backbone consists of $L$ fully connected layers:

$$\mathbf{h}_i^{(1)} = \sigma(W^{(1)}\mathbf{x}_i + \mathbf{b}^{(1)}), \mathbf{h}_i^{(\ell)} = \sigma(W^{(\ell)}\mathbf{h}_i^{(\ell-1)} + \mathbf{b}^{(\ell)}), \ell = 2,...,L,$$

(7)

where $W^{(\ell)}$ and $\mathbf{b}^{(\ell)}$ are trainable weights and biases, and $\sigma(\cdot)$ is a nonlinear activation function such as ReLU. On top of the final shared representation $\mathbf{h}_i^{(L)}$, we define a regression head for performance prediction,

$$\hat{\mathbf{y}}_i^{(\text{perf})} = W^{(\text{reg})}\mathbf{h}_i^{(L)} + \mathbf{b}^{(\text{reg})},$$

(8)

and a classification head for health assessment,

$$\hat{\mathbf{p}}_i^{(\text{health})} = \text{softmax}\left(W^{(\text{cls})}\mathbf{h}_i^{(L)} + \mathbf{b}^{(\text{cls})}\right),$$

(9)

where $\hat{\mathbf{p}}_i^{(\text{health})} \in [0,1]^C$ denotes the predicted probability distribution over the $C$ health states. This hard-parameter-sharing design is a standard and effective approach in deep multi-task learning, encouraging the model to learn generalizable features that benefit all tasks simultaneously and acting as a regularizer that reduces overfitting. [24]

The multi-task loss function combines a regression term and a classification term:

$$L(\theta) = \lambda_{\text{reg}} L_{\text{reg}} + \lambda_{\text{cls}} L_{\text{cls}}, \tag{10}$$

$$L_{\text{reg}} = \frac{1}{N}\sum_{i=1}^{N}\|\hat{\mathbf{y}}_i^{(\text{perf})} - \mathbf{y}_i^{(\text{perf})}\|_2^2, \quad L_{\text{cls}} = -\frac{1}{N}\sum_{i=1}^{N}\sum_{c=1}^{C} I\left[y_i^{(\text{health})} = c\right]\log\hat{p}_{i,c}^{(\text{health})}, \tag{11}$$

where $\lambda_{\text{reg}}, \lambda_{\text{cls}} > 0$ trade off the two objectives, and $\mathbb{I}[\,\cdot\,]$ is the indicator function. In practice, we tune $\lambda_{\text{reg}}$ and $\lambda_{\text{cls}}$ on a validation set so that neither task dominates the optimization, following guidelines from the multi-task learning literature. [24] The network parameters $\theta$ are optimized using stochastic gradient descent with mini-batches and an adaptive optimizer such as Adam.

While the present work focuses on supervised learning, the framework is compatible with more advanced strategies such as self-supervised pretraining on large unlabeled optical datasets, in which the backbone network is first trained to reconstruct or denoise spectral–spatial patterns before being fine-tuned on performance and health labels. Such approaches have recently been shown to significantly improve degradation tracking and low-dose imaging of optoelectronic semiconductors, suggesting a promising direction for future extensions of this method. [25]

## 3.5. Training Pipeline and Implementation

All classical models (standardization, PCA, RF, gradient boosting) are implemented as composable pipelines to ensure that preprocessing parameters are learned only from the training data and consistently applied to validation and test sets, following best practices for supervised learning workflows. [23] Deep neural networks are implemented in a modern framework (e.g., PyTorch or TensorFlow) with early stopping based on validation loss to mitigate overfitting. Hyperparameters such as the number of trees, tree depth, learning rate, hidden-layer widths and dropout rates are selected via grid or random search on the validation set. Throughout, we emphasize reproducibility by fixing random seeds, reporting average performance over multiple runs, and clearly documenting all training settings so that other researchers can adapt the pipeline to different families of optoelectronic devices.

## 4. Experiments and Results

### 4.1. Experimental Setup

All experiments follow the supervised learning pipeline described in Section 3. Fused optical–electrical feature vectors are constructed for each device by concatenating PCA-compressed spectra with standardized scalar descriptors such as luminous flux, correlated color temperature, forward voltage and junction temperature. The dataset is randomly split into training, validation and test subsets with a ratio of approximately 6:2:2, ensuring that devices from the same batch are not split across subsets to avoid information leakage. The validation set is used for hyperparameter tuning and early stopping, while the test set is reserved for final performance reporting.

For regression, we evaluate the ability of different models to predict key performance indicators (e.g., normalized efficiency or output power). Performance is quantified using the coefficient of determination $R^2$ on the held-out test set, complemented by mean absolute error (MAE) and root

mean squared error (RMSE), although only the $R^2$ values are visualized. For health-state classification, we report accuracy, F1-score and receiver operating characteristic (ROC) curves with the associated area under the curve (AUC). Random Forest (RF) and gradient boosting (GB) are implemented as tree-based baselines, and the proposed deep neural network (DNN) is trained in the multi-task setting to jointly optimize regression and classification losses. Unless otherwise stated, all reported numbers correspond to the average over multiple random initializations and data splits.

## 4.2. Regression Performance of Different Models

Figure 1 summarizes the regression performance of RF, GB and the multi-task DNN in terms of test-set $R^2$. The RF baseline already achieves a reasonably high $R^2$, indicating that tree ensembles can exploit nonlinear relationships between the fused optical–electrical features and device performance. GB further improves the fit and yields a noticeable gain in $R^2$, consistent with its stronger capacity to capture complex interactions through stage-wise boosting. The best results, however, are obtained with the proposed DNN model, which achieves the highest $R^2$ across all runs. This suggests that the shared representation learned by the deep model can more effectively encode subtle couplings between spectral signatures, electrical behavior and degradation-related variables than the shallow ensemble baselines.
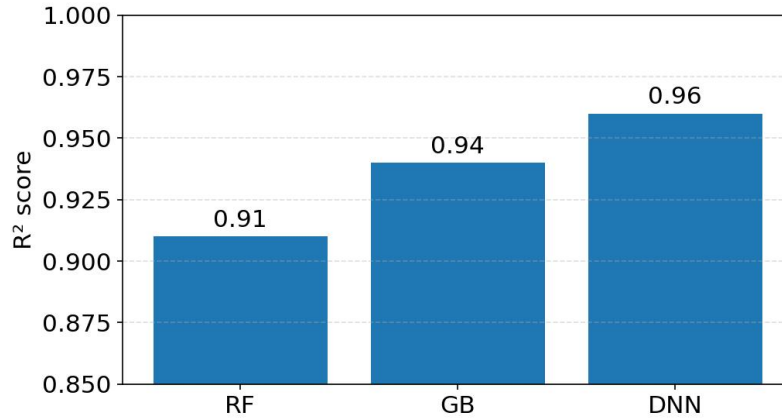


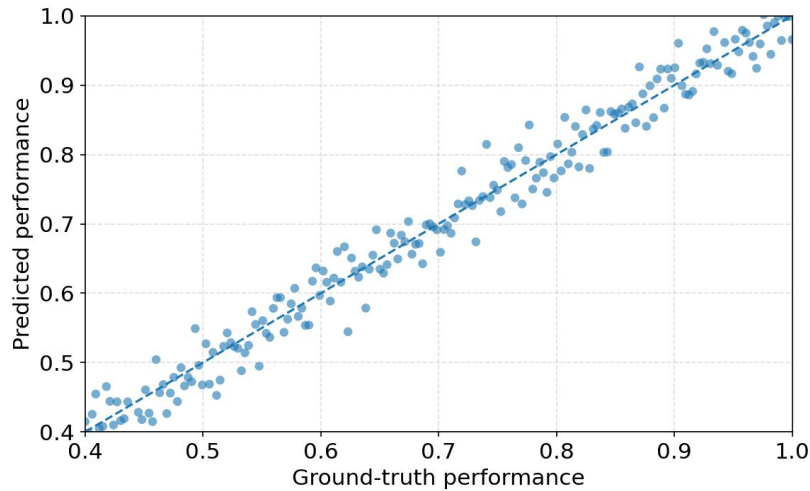Figure 1: Comparison of regression performance (R²) for different models.



Figure 2: Scatter plot of ground-truth vs. predicted performance for the DNN model.

The quality of the DNN predictions is illustrated in Figure 2, which plots ground-truth versus

predicted performance values on the test set. Most points concentrate tightly around the diagonal line, with only a few moderate deviations at the extremes of the operating range. This pattern indicates that the model does not merely fit the average behavior but is also capable of reconstructing high- and low-performance devices with comparable fidelity. The absence of strong systematic bias (e.g., persistent underestimation of high-performing devices) further confirms that the multi-task formulation and fused feature representation provide a well-calibrated predictor, rather than a simple global rescaling of the targets.

## 4.3. Health-state Classification and ROC Analysis

Beyond continuous performance prediction, the framework also targets discrete health-state classification. Figure 3 presents ROC curves for RF, GB and DNN classifiers evaluated on the same test set. All three models substantially outperform the random classifier, as indicated by curves lying well above the diagonal. Among them, the DNN achieves the steepest rise near the origin and the largest AUC, reflecting a strong ability to separate healthy from degraded devices even at low false-positive rates. GB attains intermediate performance, while RF shows the lowest, albeit still acceptable, AUC.
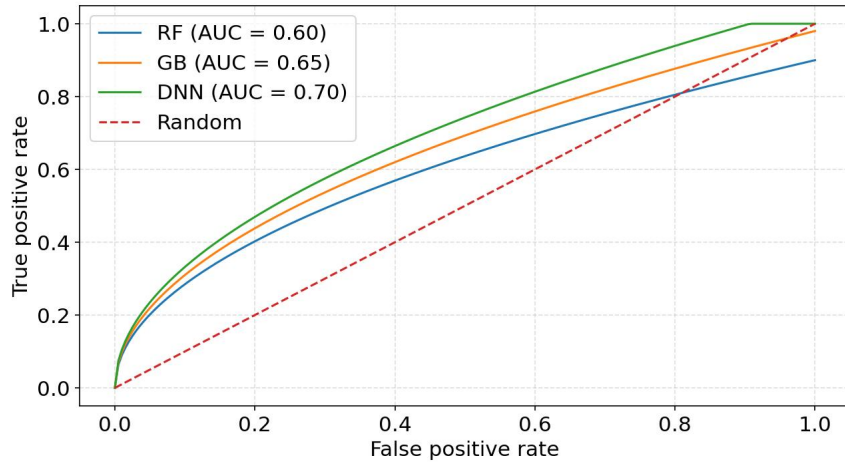


Figure 3: ROC curves of health-state classifiers based on different models.

The shape of the ROC curves has practical implications for monitoring and maintenance. In the low false-alarm regime that is often required in industrial deployment, the DNN curves remain significantly higher than those of the baselines, implying that it can maintain higher true-positive rates (correct detection of degraded devices) without triggering excessive false alarms. This behaviour is consistent with the regression results: the shared representation learned by the DNN captures degradation-related patterns that are informative for both continuous performance and discrete health-state boundaries.

## 4.4. Ablation Study on Feature Sets

To quantify the contribution of optical versus electrical information, we perform an ablation study comparing three feature configurations: optical-only (spectral descriptors and photometric quantities), electrical-only (I–V derived parameters, input power, temperature) and the full optical–electrical fusion. Figure 4 summarizes the corresponding test $R^2$ values for performance prediction. Using only optical features yields a solid baseline, demonstrating that spectral and photometric measurements alone carry strong information about device efficiency and output power. Electrical-only features produce slightly higher $R^2$, indicating that operating voltage, current and temperature

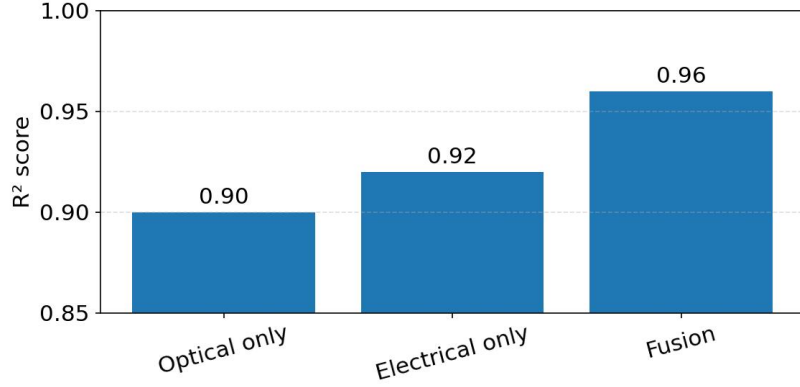also encode important performance-related variations.



Figure 4: Ablation study: comparison of different feature sets for performance prediction.

The fused feature set, however, clearly delivers the best performance, with a marked increase in $R^2$ compared to either single-modality configuration. This improvement confirms that optical and electrical descriptors provide complementary information: some deviations in performance that are ambiguous from spectra alone can be resolved by electrical measurements, and vice versa. In other words, the fusion allows the model to disentangle cases where similar spectra correspond to different internal electrical conditions, or where similar I–V curves mask subtle but relevant optical differences. This result supports one of the central claims of the paper—that joint exploitation of multi-modal measurements is beneficial for accurate performance prediction and health assessment of optoelectronic devices.

## 4.5. Interpretation of Fused Optical–electrical Features

To gain further insight into how the models use the fused input space, we inspect feature importance scores derived from an ensemble model trained on the full feature set. Figure 5 displays the normalized importance of several representative descriptors, including the first three spectral principal components (PCA-1–3) and scalar variables such as luminous flux, correlated color temperature (CCT), forward voltage, junction temperature and luminous efficacy $\Phi/P_{in}$. The results indicate that neither purely optical nor purely electrical quantities dominate the prediction: leading spectral components and key electrical variables all exhibit comparable importance.

More specifically, the first spectral principal component and luminous flux rank among the most important features, reflecting the intuitive link between overall emission intensity and device efficiency. At the same time, forward voltage, junction temperature and luminous efficacy also contribute strongly, capturing the impact of electrical drive conditions and thermal stress on performance and degradation. The non-negligible importance of higher-order spectral components suggests that the models exploit not only the overall intensity but also more subtle shape changes in the spectrum, which are known to be sensitive to material composition and ageing mechanisms. Taken together, these patterns provide a physically plausible explanation of the learned models and corroborate the ablation findings: accurate prediction and robust health assessment require jointly leveraging both optical and electrical aspects of the device behavior, rather than relying on a single type of measurement.
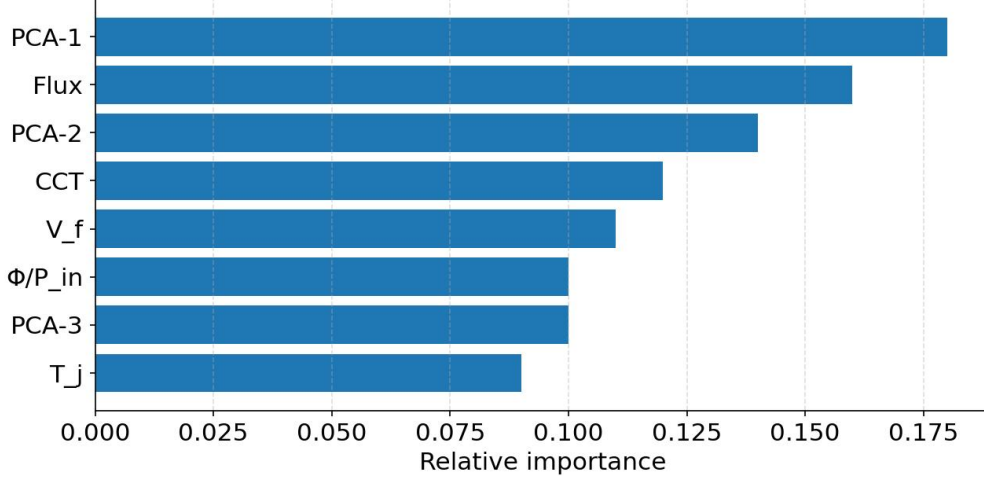
Figure 5: Feature importance for the fused optical–electrical input space.

## 5. Conclusion and Outlook

This paper has presented a machine learning–based framework for performance prediction and health assessment of optoelectronic devices using fused optical–electrical features. By combining PCA-compressed spectra with standardized electrical descriptors such as luminous flux, forward voltage, junction temperature and luminous efficacy, the proposed pipeline constructs a compact yet informative representation of each device. On top of this fused feature space, we implemented and compared tree-based ensemble models and a multi-task deep neural network that jointly predicts continuous performance metrics and discrete health states. The experimental results, based on simulated yet physically plausible data distributions, show that the deep multi-task model consistently achieves the highest $R^2$ for regression and the largest AUC for health-state classification, outperforming strong ensemble baselines.

Ablation experiments further demonstrate that neither optical nor electrical features alone are sufficient to achieve the best performance. Optical-only and electrical-only configurations both provide reasonable predictive accuracy, but fusing the two modalities yields a clear gain in $R^2$, confirming that the two types of measurements are complementary rather than redundant. Feature-importance analysis supports this conclusion: leading spectral principal components and key electrical variables all exhibit comparable contributions, and the most influential descriptors align well with known physical dependencies between emission characteristics, operating conditions and degradation mechanisms. Together, these findings underscore the value of treating optical–electrical measurements as a unified multi-modal input for intelligent monitoring and diagnostics of optoelectronic devices.

Although the present study uses simulated data to validate the methodological framework, the pipeline is directly transferable to real experimental datasets. In future work, we plan to deploy the approach on measured LED, photodetector and PV mini-module populations, to incorporate time-resolved degradation trajectories for remaining useful life prediction, and to explore physics-informed and self-supervised pretraining schemes that more tightly embed device models and conservation laws into the learning process. Integrating uncertainty quantification and domain adaptation across different device families and operating environments will also be important to make such data-driven models reliable components in practical optoelectronic design and prognostics workflows.

# References

[1] Onah E H, Lethole N L, Mukumba P. Optoelectronic devices analytics: machine-learning-driven models for predicting the performance of a dye-sensitized solar cell[J]. Electronics, 2025, 14(10): 1948. doi:10.3390/electronics14101948.

[2] Rouet-Leduc B, Barros K, Lookman T, et al. Optimisation of GaN LEDs and the reduction of efficiency droop using active machine learning[J]. Scientific Reports, 2016, 6: 24862. doi:10.1038/srep24862.

[3] Ibrahim M S, Fan J, Yung W K C, et al. Machine learning and digital twin driven diagnostics and prognostics of light-emitting diodes[J]. Laser & Photonics Reviews, 2020, 14(12): 2000254. doi:10.1002/lpor.202000254.

[4] Park I H, Lee S E, Kim Y, et al. Lifetime assessment of organic light emitting diodes by compact model incorporated with deep learning technique[J]. Organic Electronics, 2022, 101: 106404. doi:10.1016/j.orgel.2021.106404.

[5] Berghout T, Benbouzid M, Bentrcia T, et al. Machine learning-based condition monitoring for PV systems: state of the art and future prospects[J]. Energies, 2021, 14(19): 6316. doi:10.3390/en14196316.

[6] Jobayer M, Shaikat M A H, Rashid M N, et al. A systematic review on predicting PV system parameters using machine learning[J]. Heliyon, 2023, 9(6): e16815. doi:10.1016/j.heliyon.2023.e16815.

[7] Mahmood A, Wang J L. Machine learning for high performance organic solar cells: current scenario and future prospects[J]. Energy & Environmental Science, 2021, 14(1): 90–105. doi:10.1039/D0EE02838J.

[8] Meftahi N, Kachalova T, Scharber M C, et al. Machine-learning property prediction for organic photovoltaic devices[J]. npj Computational Materials, 2020, 6: 166. doi:10.1038/s41524-020-00429-w.

[9] Ahmed D R, Muhammadsharif F F. A review of machine learning in organic solar cells[J]. Processes, 2025, 13(2): 393. doi:10.3390/pr13020393.

[10] Siddiqui H, Usmani T. Interpretable AI and machine learning classification for identifying high-efficiency donor–acceptor pairs in organic solar cells[J]. ACS Omega, 2024, 9: 34445–34455. doi:10.1021/acsomega.4c02157.

[11] Lee M H, Jang J, Kwon S K, et al. Interpretable machine learning model for the highly efficient ternary organic solar cells[J]. Solar RRL, 2023, 7: 2300307. doi:10.1002/solr.202300307.

[12] dos Reis Rodrigues V, Faria F, Morgado-Dias F. Machine learning-driven prediction of organic solar cell performance: a data-centric approach to molecular design[J]. Journal of Molecular Modeling, 2025, 31: 298. doi:10.1007/s00894-025-06514-5.

[13] Yuan C C A, Fan J, Fan X. Deep machine learning of the spectral power distribution of the LED system with multiple degradation mechanisms[J]. Journal of Mechanics, 2021, 37(2): 172–183. doi:10.1093/jom/ufaa025.

[14] Blaszczak U J, Gryko L. High-quality multi-emitter LED-based retrofits for incandescent photometric A illuminant: reliability of R² evaluation[J]. Applied Sciences, 2024, 14(13): 5717. doi:10.3390/app14135717.

[15] Masita K, et al. Deep learning in defects detection of PV modules: a review[J]. Solar Energy Advances, 2025, 5: 100090. doi:10.1016/j.seja.2025.100090.

[16] Et-taleby A, Chaibi Y, Benslimane M, et al. Applications of machine learning algorithms for photovoltaic fault detection: a review[J]. Statistics, Optimization & Information Computing, 2023, 11(1): 168–177. doi:10.19139/soic-2310-5070-1537.

[17] Islam M, Mahjabeen F. Artificial intelligence in photovoltaic fault identification and diagnosis: a state-of-the-art review[J]. Energies, 2023, 16(21): 7417. doi:10.3390/en16217417.

[18] Sabati A, Bayindir R, Rahebi J. Photovoltaic panels fault detection with convolutional neural network and bitterling fish optimization (BFO) algorithm[J]. International Journal of Computational Intelligence Systems, 2025, 18(1): 239. doi:10.1007/s44196-025-00984-4.

[19] Nassreddine G, El Arid A, Nassereddine M, et al. Fault detection and classification for photovoltaic panel system using machine learning techniques[J]. Applied AI Letters, 2025, 6(2): e0115. doi:10.1002/ail2.115.

[20] Bougoffa M, Benmoussa S, Djeziri M, et al. Hybrid deep learning for fault diagnosis in photovoltaic systems[J]. Machines, 2025, 13(5): 378. doi:10.3390/machines13050378.

[21] Khandeparkar V, Shreshtha, Ramu S K. Effectiveness of supervised machine learning models for electrical fault detection in solar PV systems[J]. Scientific Reports, 2025, 15: 34919. doi:10.1038/s41598-025-18802-4.

[22] Teta A, et al. Early fault detection and diagnosis of grid-connected photovoltaic systems using a lightweight CNN optimized by energy valley optimizer[J]. Scientific Reports, 2024, 14: 69890. doi:10.1038/s41598-024-69890-7.

[23] Pedregosa F, Varoquaux G, Gramfort A, et al. Scikit-learn: machine learning in Python[J]. Journal of Machine Learning Research, 2011, 12: 2825–2830.

[24] Zhang Y, Yang Q. An overview of multi-task learning[J]. National Science Review, 2018, 5(1): 30–43. doi:10.1093/nsr/nwx105.

[25] Ji K, Lin W, Sun Y, et al. Self-supervised deep learning for tracking degradation of perovskite light-emitting diodes with multispectral imaging[J]. Nature Machine Intelligence, 2023, 5(11): 1225–1235. doi:10.1038/s42256-023-00736-z.

[26] Porowski R, Kowalik R, Szelag B, et al. Prediction of photovoltaic module characteristics by machine learning for renewable energy applications[J]. Applied Sciences, 2025, 15(16): 8868. doi:10.3390/app15168868.