

Machine Learning–Based Detection and Classification of Overlapping Fiber Anomalies

Original

Machine Learning–Based Detection and Classification of Overlapping Fiber Anomalies / Malik, Gulmina; Masood, Muhammad Umar; Ali, Ahtisham; Cheruvakkadu Mohamed, Mashboob; Straullu, Stefano; Nespola, Antonino; Kishore Bhyri, Sai; Napoli, Antonio; Ao Pedro, Jo; Maria Galimberti, Gabriele; Wakim, Walid; Curri, Vittorio. - In: IEEE PHOTONICS TECHNOLOGY LETTERS. - ISSN 1041-1135. - (2026). [10.1109/LPT.2026.3678082]

Availability:

This version is available at: 11583/3009472 since: 2026-04-01T09:00:13Z

Publisher:

IEEE

Published

DOI:10.1109/LPT.2026.3678082












Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)

Machine Learning–Based Detection and Classification of Overlapping Fiber Anomalies

Gulmina Malik , Muhammad Umar Masood , Ahtisham Ali , Mashboob Cheruvakkadu Mohamed 
Stefano Straullu , Antonino Nespola , Sai Kishore Bhyri , Antonio Napoli , João Pedro 
Gabriele Maria Galimberti , Walid Wakim and Vittorio Curri 

Abstract—This work presents a machine-learning framework for detecting and classifying independent and overlapping anomaly signatures in optical fiber infrastructure using state-of-polarization (SOP) dynamics. The framework exploits temporal variations of the Stokes parameters to characterize bending, tapping, shaking, and their simultaneous combinations. An experimentally collected dataset is used to evaluate multiple machine-learning models under realistic noise conditions representative of practical fiber monitoring environments. The results demonstrate that SOP-based features enable accurate identification of complex fiber disturbances, strengthening physical-layer monitoring and supporting the secure and reliable operation of next-generation optical networks.

Index Terms—Machine learning, state-of-polarization, optical fiber, fiber anomalies, overlapping signatures.

I. INTRODUCTION

FIBER optic networks are the backbone of global telecommunications, allowing for high-speed data transmission with minimal delay [1]. However, as these networks develop, they become more vulnerable to leakage of any sensitive and confidential data, especially overlapping anomalies like bends, breakage, and splice losses. Overlapping disturbances hamper identification, requiring more complex detection techniques [2]. Optical fibers are susceptible to unauthorized eavesdropping attempts, malicious physical attacks, fiber cuts and environmental vibrations [3], [4]. Such attacks can result in severe network disruptions and compromised confidentiality, so it is crucial to detect them earlier on. Traditional approaches, such as the use of optical time-domain reflectometry (OTDR) and threshold-based monitoring, rely on established heuristics and struggle to distinguish overlapping anomalies, especially when numerous fault signs interact [2].

Recent advances in machine learning (ML) have allowed for more advanced anomaly detection in optical networks. In [5], a vision transformer-based model was presented to recognize and locate simultaneous anomaly occurrences. However, such methods typically depend on spectrogram-based representations and deep architectures, introducing substantial

preprocessing overhead and inference latency. As a result, they are often unsuitable for real-time, low latency resource-constrained monitoring systems, deployed at the network edge. These problems highlight the necessity for intelligent real-time anomaly detection framework that can take preventive measures without jeopardizing network integrity or causing an economic disruption. It is critical to monitor the metropolitan environment in order to provide early warnings for scenarios that might compromise the optical networks' health. As an early warning, the architecture should be able to classify the type of anomaly and track it to mitigate its impact, as much as possible, by implementing preventive steps [6]. While deep learning models are capable of learning complex nonlinear relationships, our results show that they do not inherently provide lower per-sample inference latency than traditional machine-learning approaches, particularly in resource-constrained SOP monitoring systems.

State-of-polarization (SOP) monitoring is a promising ML-assisted technique for detecting mechanical and geophysical disturbances in optical fibers. [7]. The SOP represents the orientation of the optical field, and its variations provide a sensitive indicator of external perturbations along the fiber. The SOP angular speed (SOPAS), captures rapid polarization dynamics on the Poincaré sphere. This makes it particularly effective for identifying fine-grained events such as bending, tapping, and different forms of shaking, bending, vertical oscillation of fiber and tapping (light mechanical hit) [8]. However, current SOP-based techniques are insufficiently detailed in their classification of overlap irregularities, and they frequently fail to distinguish between hostile intrusions (such as eavesdropping) and harmless natural vibrations. While previous studies, such as [9], have successfully utilized ML to classify overlapping events, several critical gaps remain for practical deployment.

Detection and classification of overlapping anomalous events is a complex challenge. We extend our work in [10], by incorporating a wider range of events and testing multiple models, we can find the best methods in terms of performance and computational needs, enabling a lightweight pipeline suitable for real-time edge deployment in real-world metro transport network. The results show that XGBoost provides the optimum balance of accuracy, inference, and computing resources, making it appropriate for low-latency, per-sample inference in metro-scale SOP monitoring systems.

Gulmina Malik (gulmina.malik@polito.it), Muhammed Umar Masood, Ahtisham Ali, Mashboob Cheruvakkadu Mohamed and Vittorio Curri are with Politecnico di Torino, Italy. Stefano Straullu is with Links Foundation, Italy and Sai Kishore Bhyri, Antonio Napoli, João Pedro, Gabriele Maria Galimberti and Walid Wakim are with Nokia - Optical Networks. João Pedro is also with Instituto de Telecomunicações, Instituto Superior Técnico, Portugal. This work has received funding from the project PNRR-NGEU (MUR-DM117/2023), and from the EU's Horizon Europe research and innovation program under GA No. 101092766 (ALLEGRO Project), and EWOC GA No. 101073265.

Manuscript received Nov xx, 2025; revised Feb xx, 2026.

0000-0000/00\$00.00 © 2025 IEEE

TABLE I
SIGNATURES COLLECTED FOR ANOMALY DETECTION

Events	Justification
<i>bl</i>	Baseline; normal operating conditions
<i>b</i>	Bending;
<i>tp</i>	Tapping or Small hit
<i>ss</i>	Slow shaking
<i>fs</i>	Fast shaking
<i>ss_tp</i>	Shaking + Tapping
<i>ss_b</i>	Shaking + Bending

TABLE II
ACCURACY AND INFERENCE TIME ANALYSIS OF ML MODELS

Model	Training Accuracy(%)	Test Accuracy(%)	Inference Time(μ s)
RF	100	99.9	12713.14 \pm 0.0552
XGBoost	100	99.9	122.477 \pm 0.0344
LR	99.76	7.05	25.379 \pm 0.0022
DT	100	96.2	35.60 \pm 0.0023
GB	100	100	399.665 \pm 0.3110
MLP	99.99	40.9	74.324 \pm 0.0489
k-NN	99.8	69.08	3205.096 \pm 6.0845
SVM	99.76	29.7	115.534 \pm 3.4872

II. DATA COLLECTION AND ANOMALY DETECTION SETUP

The current study entails gathering data on seven separate signatures, including overlapping signatures, in an installed fiber link, in the LINKS Foundation [11] lab. We used the unique polarization fingerprint of each event as a key metric to distinguish between concurrent events that occur at the same time. The experimental testbed consists of two G.652 standard single-mode fibers (SSMF) that stretch 13 kms. A continuous wave (CW) of 1530 nm light pulses is injected into the sensing fiber. A Novoptel PM1000 polarimeter is linked at the receiving end to trace Stokes parameters $\{S_0, S_1, S_2, S_3\}$ to generate and monitor the polarization signatures on the Poincaré sphere. These Stokes parameters define the polarization state of light, and since S_0 represents the total optical intensity and carries no polarization information, only $S_1, S_2,$ and S_3 were considered in the further analysis.

Building on our work in [10], we enriched our dataset by incorporating more anomalous events, emulating real-world perturbations, as seen in Table I. “*bl*” represents relaxed/stable or baseline fiber with no disturbance. For the other events, the Arduino-controlled robotic arm was programmed to generate precise and reproducible perturbations, such as “*tp*” with 1 tap per second and “*ss, fs*” referred to as shaking with varying frequencies, simulating fiber vibrations, as detailed in [8].

Specifically, shaking was induced by driving the arm in a vertical up-and-down motion at a 90° angle of deviation relative to the fiber axis. Slow shaking (*ss*) is done at the frequency range of 1-5 hz and fast shaking (*fs*) is done in the range of 5-10 hz, to mimic environmental vibrations such as wind-induced oscillations and nearby vehicular or construction activities. In addition, we employed an optical fiber identifier (OFI) to classify bending “*b*” or eavesdropping, with a bend diameter of 0.25 mm. We clamped for 5 seconds and then released, causing light leakage. This test was conducted 25 times to guarantee robustness. To synchronize overlapping events, we coordinated the robotic arm’s fiber shaking action with bending, “*ss_b*” and fiber hitting, “*ss_tp*”. While traditional methods often combine overlapping events into a single class, we specifically regard overlapping events as unique, separate classes, resulting in more granular diagnostic capabilities. Following the acquisition of this data, we collected 324,117 samples for all the events. The dataset was unbalanced due to more repetitive counts of some events, with the *bl* class greatly out-representing the other anomalous events. To solve this, we performed a few preprocessing steps before feeding the dataset into our framework. We first normalized the raw data samples to a uniform scale to ensure consistency and uniformity, and then filtered out insignificant transient values.

Furthermore, to reduce class imbalance, we used synthetic minority over sampling technique (SMOTE) and class weight adjustment techniques. These techniques helps to keep the models sensitive to rare anomalous signatures.

III. MACHINE LEARNING MODELS

A. Machine Learning Architecture

Our proposed architecture for the classification of overlapping malicious events integrates the detection of distinct anomaly signatures using polarization fingerprints with advanced ML techniques, using our lab’s experimental testbed. To achieve granular detection, we used a multi-class output classifier technique, where overlapping events are handled as distinct, separate events rather than a single aggregate class. Our analysis included the following classifiers from the scikit-learn library: XGBoost, random forest (RF), gradient boosted trees, decision tree (DT), logistic regression (LR), support vector machine (SVM), k-nearest neighbor (k-NN), and neural networks: multi-layer perceptron (MLP). The dataset was divided into 80% for training and 20% for testing of the models. The models were evaluated on a Apple M4 processor with 24 GB of RAM. We only trained the models with the top features selected, based on our feature importance analysis, discussed in detail in Sec. III-B and eliminated less relevant features to speed up the training process and increase the models’ sensitivity. During the testing phase, we used a synthetic multi-component noise model to evaluate robustness under actual operating settings, simulating environmental micro-vibrations that are common in practical deployments. The model incorporates low-frequency traffic disruptions, wind-induced pink noise, and narrow-band equipment vibrations with minimal frequency modulation.

Each model was evaluated based on key performance metrics such as accuracy, F-1 score, and computational load, across a dedicated testing dataset, which is essential for real-time low latency edge applications.

B. Feature Engineering

The input data consists of time-series measurements of the Stokes parameters $\{S_1, S_2,$ and $S_3\}$. Because of its fine-grained temporal resolution, the framework can detect or capture even the smallest polarization transients that could be signs of external invasions like eavesdropping or vibrational disturbances. The models are trained first on clean data including all events described in Table I. During the training process, we further apply extensive feature engineering to add more characteristics and assure our models’ robustness, generalization, and scalability. The raw polarization data are

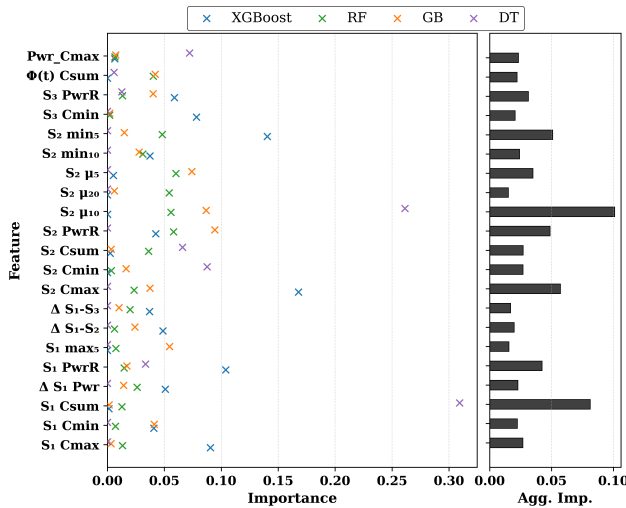


Fig. 1. Feature Importance for the Top Models

translated into a more complex feature representation by including power and SOPAS, which captures the strength and rate of polarization shift. Rolling window statistics for each Stokes parameter were computed, which included mean, standard deviation, min, max, and range, over window sizes of [5,10,20], to identify additional patterns for each window feature. To quantify the distribution and shape of the data, we calculated some additional advanced statistical features, including skewness, kurtosis, z-scores, and percentiles. We also incorporated differentials of all the Stokes parameters and SOPAS up to third order to capture the short-term dependencies, and cumulative sum and min/max to capture long-term trends. These predefined features provide a complete representation of the data, allowing the ML algorithms to find and categorize overlapping irregularities in the fiber optic infrastructure.

Fig. 1 depicts the most impactful and beneficial features for our training, by plotting each feature along y-axis and its relevance score for various models at x-axis. We applied feature engineering analysis on all the models but selected the top four models (Table II) based on their performance and displayed their feature importance score and assigned each a different color. The scatter plot indicates the importance score for each model, while the bar chart on right, displays the overall relevance of each feature across all the models. The plot shows that “ $S_2\text{-}\mu_{10}$ ” and “ $S_1\text{ Csum}$ ” have the highest feature importance in all the models. This analysis assists in providing the foundation for selecting the most relevant features to consider while tweaking the models in our later analysis in Sec. IV, for further optimization of input vector, reducing computational redundancy, saving time and resources.

IV. PERFORMANCE COMPARISON OF ML MODELS

In this section, we evaluate the efficacy of various models for classifying and detecting fiber anomalies in optical networks. We analyzed a number of ML models and compared their performance using various metrics, such as accuracy, F-1 score, computational power, training time, inference time, and FPR. We investigated the temporal fluctuation of Stokes parameters, which preserves the different polarization finger-

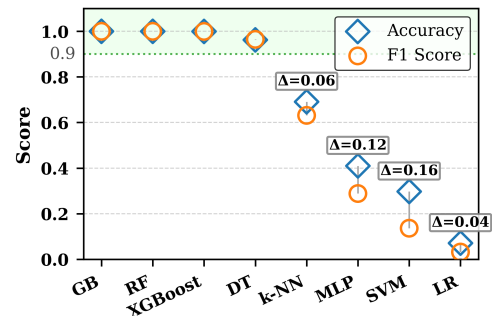


Fig. 2. Accuracy and F-1 score of various models

prints required to detect simultaneous disturbances, including the overlapping events, and demonstrated overall accuracy across those complex signatures. Our detailed analysis emphasizes the significance of choosing models that achieve a balanced performance in terms of accuracy and F-1 score and computational load for real-network anomaly detection tasks.

Table II summarizes the performance of various ML models, emphasizing the requirement for selection of efficient model in real-world deployment. It should be noted that to evaluate the feasibility of real-time detection, we measured the per-sample inference latency. Real-time refers to low-latency, per-sample inference within the monitoring pipeline, where each SOP measurement is classified immediately upon acquisition, without batch processing or offline aggregation. RF, XGBoost and GB, perform exceptionally well with near-perfect accuracy, exceeding 99.9%. Conversely, some models prove to be inefficient for the task, due to high inference time. Importantly, this table illustrates the trade-off between efficiency and inference time, highlighting the importance of efficient model selection in real-world deployment, as per requirements. Therefore, GB and XGBoost both prove to be the best options by exhibiting the best balance between high accuracy and quick inference. On the other hand, the other models show performance constraints on either metric.

Another representation of accuracy and F-1 score of these models is illustrated in Fig. 2. F-1 score depicts the harmonic mean of precision and recall and is particularly useful when dealing with unbalanced datasets. As shown, all the three models – RF, XGboost and GB – perform remarkably well, reaching the threshold of one for both accuracy and F-1 score. Following these, DT also performed well with a moderate variance in accuracy and F-1 score from the other models, indicating some limitations in handling overlapping anomalies. For overlapping events specifically, XGBoost further outperforms, achieving an F-1 score of 0.991 for “ ss_{tp} ” and 0.996 for “ ss_b ”, highlighting its strong recall-oriented detection capability in simultaneous anomalous scenarios. As also visible in Fig. 2, SVM, LR, k-NN, MLP under-perform when dealing with complex and noisy data. Fig. 3 depicts the Receiver Operating Characteristic (ROC) curves for all the models we have used, illustrating their True Positive Rate (TPR) versus False Positive Rate (FPR). The Area Under the Curve (AUC) values determines the model’s ability in distinguishing between the classes. ROC provides insight into the model’s discriminative power for each event type, regardless of a fixed classification threshold. Noticeably, the

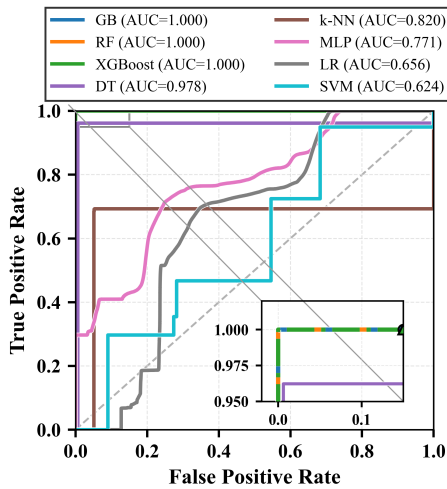


Fig. 3. ROC curves for various ML Models

tree based models outperformed other models such as k-NN (AUC = 0.820), MLP (AUC = 0.771), LR (AUC = 0.656), and SVM (AUC = 0.624), with respect to distinguishing between true and false positives. The inset shows a zoomed-in picture to highlight the differences in performance across various models with lower false positive rates.

The radar plot in Fig. 4 provides a comparative view of the accuracy–efficiency trade-offs across the evaluated machine learning models. The green dashed curve illustrates the model accuracy (normalized between 0 and 1), where larger radial values indicate superior detection performance. The three solid curves represent the Computational Load (CL) under different weighting configurations, each aligned with a distinct deployment scenario. The weighted CL metric for each scenario k is given by:

$$CL_k = \sum_j (w_{k,j} \tilde{x}_j), \quad k \in \{1, 2, 3\}.$$

where $w_{k,j}$ denotes the weight assigned to the j -th normalized computational feature \tilde{x}_j . The feature set includes inference time, training time, average and peak CPU utilization, and memory usage. The three weighting schemes capture different system priorities: (i) w_1 emphasizes inference and training time, reflecting latency-sensitive environments; (ii) w_2 provides a balanced contribution from all features; and (iii) w_3 increases the influence of CPU and memory usage, representing resource-constrained platforms such as embedded or edge-based fiber-monitoring systems. As can be clearly seen, tree-based ensemble methods offer the best trade-off, featuring very high accuracy with the use of a small amount of processing power, thus being suitable for real-time applications. On the other hand, k-NN and SVM are computationally expensive and offer lower accuracy, which makes them unsuitable for high speed fiber monitoring systems. Overall, this plot clearly observes that both GB and XGBoost are the most effective models for anomaly detection in real-time applications in fiber-optic networks. While both exhibit high accuracy, XGBoost offers a distinct advantage in inference speed and computational efficiency. This superior balance makes it the most

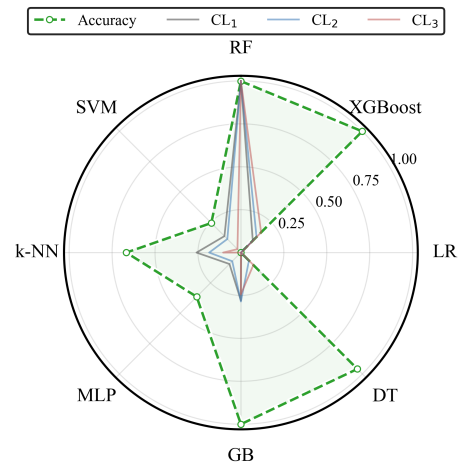


Fig. 4. Radar plot illustrating accuracy-efficiency trade-offs of various ML models

suitable candidate for the real-time fiber sensing applications in deployed resource-constrained infrastructure.

V. CONCLUSION AND FUTURE WORK

This paper demonstrates that polarization-based features, particularly those derived from the temporal evolution of Stokes parameters, provide a strong basis for detecting both normal and overlapping anomalies in optical fiber infrastructure. Among the evaluated models, tree-based ensemble methods, especially XGBoost, achieve the best balance of accuracy, robustness under noise, and fast inference, enabling reliable real-time identification of both individual and overlapping events. The approach presented here establishes a strong baseline for polarization-based anomaly detection, and future studies may extend this foundation to testing the model in an operational metro transport network.

REFERENCES

- [1] L. X., “Evolution of fiber-optic transmission and networking toward the 5g era,” *iScience*, vol. 22, pp. 489–506, 2019.
- [2] K. Abdelli *et al.*, “Machine-learning-based anomaly detection in optical fiber monitoring,” *Journal of optical communications and networking*, vol. 14, no. 5, pp. 365–375, 2022.
- [3] C. Natalino *et al.*, “Experimental study of machine-learning-based detection and identification of physical-layer attacks in optical networks,” *Journal of Lightwave Technology*, vol. 37, no. 16, pp. 4173–4182, 2019.
- [4] Y. Li *et al.*, “Fiber eavesdropping detection and location in optical communication system,” *Photonics*, vol. 12, no. 5, 2025. [Online]. Available: <https://www.mdpi.com/2304-6732/12/5/501>
- [5] K. Abdelli *et al.*, “Vision transformers for anomaly classification and localization in optical networks using sop spectrograms,” *Journal of Lightwave Technology*, vol. 43, no. 4, pp. 1902–1914, 2025.
- [6] G. Malik *et al.*, “Resilient anomaly detection in fiber-optic networks: A machine learning framework for multi-threat identification using state-of-polarization monitoring,” *AI*, vol. 6, no. 7, 2025. [Online]. Available: <https://www.mdpi.com/2673-2688/6/7/131>
- [7] A. Mecozzi *et al.*, “Polarization sensing using submarine optical cables,” *Optica*, vol. 8, no. 6, pp. 788–795, Jun 2021. [Online]. Available: <https://opg.optica.org/optica/abstract.cfm?URI=optica-8-6-788>
- [8] G. Malik *et al.*, “Machine learning for predictive multi-event detection in fiber optic systems,” in *International Conference on Machine Learning and Communication Networks (ICMLCN)*, 2025.
- [9] L. Sadighi *et al.*, “Machine learning-based polarization signature analysis for detection and categorization of eavesdropping and harmful events,” in *2024 Optical Fiber Communications Conference and Exhibition (OFC)*. IEEE, 2024, pp. 1–3.
- [10] G. Malik *et al.*, “Intelligent detection of overlapping fiber anomalies in optical networks using machine learning,” in *2025 IEEE Photonics Society Summer Topicals Meeting Series (SUM)*, 2025, pp. 1–2.
- [11] [Online]. Available: <https://linksfoundation.com/en/>