Online Dynamic Cyber-Attack Diagnosis in Power Electronics Systems Based on Few-Shot Learning

Qi Li¹, Jinan Zhang¹, Jin Ye¹, Liang Zhao², Tianqi Hong¹, Jamie Lian⁵, Beshoy Morkos¹, Hongyue Sun¹, Feraidoon Zahiri³, Chris Farnell ⁴, Alan Mantooth⁴, and WenZhan Song¹

¹Center for Cyber-Physical Systems, University of Georgia, Athens, Georgia 30602

²College of Computing and Software Engineering, Kennesaw State University, Marietta, GA 30060

³402 CMXG/MXDEO, Robins Air Force Base, Warner Robins, GA 31098

⁴Department of Electrical Engineering, University Arkansas, Fayetteville, AR 72701

¹{qi.li, jinan.zhang, jin.ye, tianqi.hong, bmorkos, hongyuesun, wsong}@uga.edu,

²lzhao10@kennesaw.edu, ⁴cfarnell, mantooth@uark.edu, ³feraidoon.zahiri@us.af.mil

⁵lianj@ornl.gov

Abstract— With increasing exposure to software-based sensing and control, power electronics systems are facing higher risks of cyber-physical attacks. To ensure system stability and minimize potential economic losses, it is critical to monitor the operating states and detect those attacks at the early stage. However, anomaly detection and diagnosis of attacks are still challenging, especially when labeled anomaly data is difficult or even infeasible to obtain. To overcome this problem, we propose a Few-Shot Learning (FSL) based approach for cyber-attack diagnosis leveraging the waveform data. To the best of our knowledge, this work is the first attempt at leveraging FSL for cyber-attack diagnosis in power electronics systems. Extensive experimental results demonstrate that our proposed approach can achieve comparable diagnosis accuracy with the state-of-the-art data-driven methods using less than 0.04% of the training samples.

Index Terms—Few-shot learning, power system, cyber-attack, attack diagnosis, siamese neural network, deep learning

I. INTRODUCTION

With the booming popularity of the Internet in modern society, a huge amount of devices are connected through networks, and the security of cyberspace has drawn great attention at present, including the area of power systems. A variety of datadriven methods have been widely adopted for attack and event detection in power electronics-based smart grids, including rule-based data-driven analytics [1] and signal-property-based approach [2]. In particular, there is proven success in adopting machine learning and deep learning methods for enhanced performance [3], [4]. For instance, the multilayer long shortterm memory networks have been used to leverage time-series electric waveform data from current and voltage sensors for attack detection [5]. Zhao et al [6] proposed a novel privacypreserving decentralized detection framework incorporating federated learning (FL) that enables collaboratively training across devices without sharing raw data.

However, the performance of those methods is highly dependent on the data domain. First, a huge volume of the training dataset is required to feed the model that may not be available in practical applications, especially for anomaly or attacked cases. Second, for model adaptation with new data, retraining the entire model is needed, which is time-consuming

and not suitable for time-sensitive and resource-constrained applications. Moreover, in the scenario of cyber-attack detection and diagnosis in the power system, labeled attacked data is usually difficult or even infeasible to obtain, and it is common that new types of attacks occur constantly due to the dynamics of the cyberspace environment. For instance, zero-day attacks launch on the day when a vulnerability is discovered [7]. The versatility of the environment makes it challenging to guarantee cyber-attack diagnosis performance as the data-driven model needs to adapt quickly to new attacks.

To address the issues above, we aim to close this gap by developing a Few-Shot Learning (FSL) based paradigm that can accurately detect cyber-attacks using only a few labeled samples in power electronics systems, and quickly adapt to new unseen attacks. In the literature, FSL is mainly focused on supervised learning problems such as few-shot classification for image classification [8], and object recognition [9] when very limited labeled data is available. However, there is little research using FSL on power electronics system applications. Based on our previous work [5], [10], [11] that successfully used data-driven approaches for cyber-attack detection in power systems, we further adopt FSL for attack diagnosis in power electronics systems in this work.

The contributions of this work are summarized as follows:

- This is the first work of data-driven methods for attack detection in power electronics that considers the constraint of limited labeled data and addresses the newly identified category of unseen attacked data.
- 2) We developed a novel FSL-based diagnosis framework that requires only a few labeled samples and can quickly adapt to a dynamic environment with unseen attacks.
- 3) Extensive experiments are conducted to validate the suitability and efficiency of the proposed method in diagnosing cyber-attacks in solar farm case studies.

II. SYSTEM AND ATTACK MODEL

A. System Modeling

To explore the cyber-attack in the power system, we specifically simulated a solar farm system and proposed the cyber-

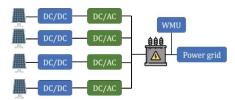


Fig. 1. Schematic diagram of the PV farm Simulink model. An attacker is able to launch cyber-attacks between the controller and DC/DC converter or DC/AC inverter. The WMU is installed between the transformer and the power grid, which monitors the raw waveform data.

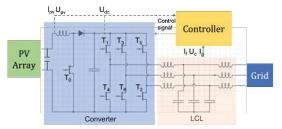


Fig. 2. Two-stage two-level PV converter circuit.

attack model to generate the dataset based on the same setting as our previous work [5]. Fig. 1 shows the schematic diagram of our solar farm Simulink model. Fig. 2 illustrates one converter circuit as an example. The main power grid is modeled as an ideal voltage source, and the load is linear. One rate voltage of 260V/25kV, 400kVA, transformer connects the PV farm to the power grid. Note that we are mainly focusing on the scenario in which a cyber-attack occurs on the DC/DC converters and DC/AC inverters, which would bring in unusual harmonics and then affect the power quality in the power systems. In addition, the waveform measurement unit (WMU) is installed between the transformer and the power grid. It monitors changes in the characteristics of the power system when the DC/AC inverter and DC/DC converter are under attack.

B. Attack Model

We proposed the following cyber-attack model:

$$Y_F(t) = \alpha Y_0(t - t_{delay}) + \beta, \tag{1}$$

where Y_F is the manipulated data vector that is the input of the controller; Y_0 is the original measurement; α is a multiplicative factor matrix that defines the weight of the attack vector; β is a multiplicative factor that defines the weight of the real vector; t_{delay} is the delay time injection. Our cyber-attack models effectively reflect real-world cyber threats. Validation through hardware experiments and existing literature ensures that our models align with actual cyberattacks observed in similar systems. Fig 3 shows an example of a DC/AC inverter attack from our PV system. Nevertheless, our simulation results demonstrate that cyber-attacks could also result in subtle distortions, which makes detection and diagnosis tougher. To obtain an accurate model, we have simulated the following three cases: 1) Normal condition. 2) False data injection attack on the PV converter controller. 3) Delay attack on PV converter. With the variation of irradiation levels and attack parameters, we end up with 9 cases in total. The details will be discussed in IV-A.

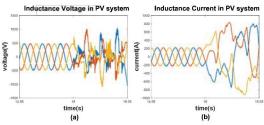


Fig. 3. Electric waveforms (voltage and current) simulations of a DC/AC inverter attack example from our PV system.

III. ALGORITHM DESIGN

A. Online Dynamic Attack Diagnosis

To address the concerns outlined in Section I and tackle the challenge of learning from a limited number of samples, we propose an online dynamic attack diagnosis approach based on FSL [12]. The fundamental concept behind FSL is to train a model that can learn the similarity between two input samples.

During the model prediction stage, the FSL model is presented with a classification task denoted as T, along with a dataset $D = \{D_{query}, D_{support}\}$. The support set $D_{support}$ comprises candidate data samples and their corresponding labels, represented as $\{(x_i, y_i)\}_{i=1}^{I}$, with I indicating the small size of the support set (usually fewer than 10). The query set, denoted by D_{query} , comprises data instances that need to be classified. These are symbolized as $\{(x_j)\}_{j=1}^{J}$, where J signifies the total number of queries. The FSL model computes the similarity between the data samples in the query set and those in the support set. Subsequently, it assigns the query data to the class of the most similar candidate within the support set. Conventionally, one considers the N-way-K-shot setting [12], in which $D_{support}$ contains I = KN examples from N classes each with K candidate data samples.

During the training stage, the parameter θ_0 of the model is optimized to minimize the error across training data, where the error is defined by the loss function such as contrastive loss, and triplet loss [13]:

$$\theta_0 = \arg\min_{\theta_0} \sum_{((x_i, y_i), (x_j, y_j)) \in D_{train}} \ell(h(x_i, x_j; \theta_0), \hat{y}), \quad (2)$$

where h denotes the FSL model, D_{train} denotes the training set, \hat{y} represents whether y_i and y_j belong to the same category, and ℓ is the loss function. The loss function aims to maximize the proximity of data instances belonging to the same class in the feature space, while simultaneously encouraging a clear separation between instances from different classes. It is noteworthy that the FSL model has never been exposed to either the support set or the query set during training. Additionally, the model has not encountered any specific classes represented by the query set. This highlights a notable advantage of the FSL model over conventional learning approaches: the ability to classify new data classes with limited samples, without the need for retraining.

Fig 4 illustrates the workflow of our proposed online cyberattack diagnosis framework. Two FSL models are trained: one focuses on detection, while the other is dedicated to diagnosis.

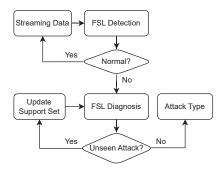


Fig. 4. Workflow of online dynamic attack diagnosis.

The FSL detection model serves as an online attack detector to analyze the streaming data. Upon identifying an anomalous data segment, the detection model promptly transfers it to the FSL diagnosis model. The latter then diagnoses the attack type. Explicitly, during diagnosis, each streaming data segment is fed into the FSL model and compared with candidates in the support set. If all comparison results exceed a predefined threshold, the segment is identified and labeled as a new attack type, prompting an update to the support set with this new category. We developed the online diagnosis algorithm as Algorithm 1.

Algorithm 1 Algorithm for Online Cyber-Attack Diagnosis

```
1: Function DIAGNOSEATTACK
   Input:
     th: a predefined threshold to determine attack categories.
      sSet: the support set.
2: sScore \leftarrow \emptyset
                              3: while true do
       dS \leftarrow \text{getNextDataSegment}()
4:
       isAbnormal \leftarrow DetectionFSL(dS)
5:
       if not is Abnormal then
6:
           continue
7:
       end if
8:
       for c \in sSet do
9:
           sScore[c] \leftarrow calculateSimilarity(dS, c)
10:
       end for
11:
12:
       if min(sScore) < th then
           return arg min(sScore)
13:
14:
       else
           newCategory \leftarrow createNewCategory(dS)
15:
           addCategoryToSupportSet(newCategory, sSet)
16:
           return newCategory
17:
       end if
18
19: end while
```

B. Transferable Siamese Neural Network

The Siamese Neural Network (SNN) [14] is adopted to implement the FSL model. Fig 5a shows the conventional structure of SNN. Pairs of input are fed into the model that shares the weights to train the feature embedding. All the parameters in the dashed box are trainable. Once the network is well-trained with the capability to capture the similarity between two inputs, it will be able to predict new classes with

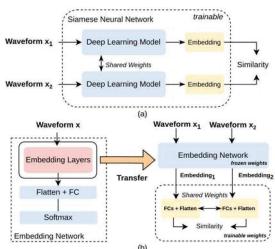


Fig. 5. (a) conventional SNN architecture for implementing FSL, where all parameters in the dotted box are trainable. (b) Proposed TSNN architecture for implementing FSL. The embedding network acts as a frozen weights backbone, which is transferred from the network for classification.

only a few samples stored in the support set. However, SNN is trained on pairs of inputs, making the training more complex and potentially time-consuming compared to traditional neural networks, and the performance heavily relies on the chosen similarity metric.

- 1) Network Architecture: To enhance the performance and convergence efficiency, we modified the SNN and proposed a Transferable Siamese Neural Network (TSNN) structure shown as Fig 5b by utilizing Transfer Learning. The training stages of our proposed TSNN architecture exhibit enhanced convergence efficiency and faster training speed per epoch due to the reduced number of trainable parameters resulting from the application of Transfer Learning. Specifically, the training of TSNN will be divided into two steps: (a) Embedding network training. This step focuses on training the embedding network, where a deep learning model is employed to train a classifier using the available training data. This training process continues until the classification results reach an acceptable level, signifying that the model has acquired the ability to effectively represent the distinct classes present in the training data. Then the model of the embedding network is saved for the next step. In our case, we use a Convolutional Neural Network (CNN) with residual block thanks to its remarkable classification performance. (b) Fine-tuning. In this step, the pre-trained embedding network serves as a frozen backbone for the SNN. The embedding output is connected to trainable fully-connected layers.
- 2) Loss Function Design: We employ a proposed weighted loss function 3 to optimize the training process:

$$L(x_1^i, x_2^i) = (1 - \alpha)L_1(x_1^i) + \alpha L_2(x_1^i, x_2^i),$$
(3)

where L_1 is the cross-entropy loss, i indexes the i-th minibatch of training data, α is the weight between 0 to 1, and L_2 is the contrastive loss function defined as:

$$L_2 = (1 - Y) * ||x_1 - x_2||^2 + Y * \max(0, m - ||x_1 - x_2||^2),$$
(4)

 $\label{table I} \textbf{TABLE I} \\ \textbf{F1-score for the 3, 7, 10-shots attack detection}$

Support set	3-shots	7-shots	10-shots
attack 1	0.893	0.895	0.898
attack 2	0.901	0.912	0.920
attack 3	0.881	0.891	0.903
attack 4	0.875	0.881	0.897
attack 5	0.912	0.902	0.914
attack 6	0.852	0.869	0.885

where Y=1 whenever x_1 and x_2 are from the same class and Y=0 otherwise, and m, the margin value, determines when the dissimilar classes are far enough apart. The initial weights are drawn from a standard normal distribution, and biases were drawn from a Gaussian with a mean of 0.5 and a fixed variance of 10^{-2} .

IV. EXPERIMENTAL EVALUATION

In this section, we conducted experiments on two study cases to evaluate our method's performance. First, we analyze its performance based on the model in Section II. Then, we explore its adaptability in a more challenging case. Both cases are detailed in Section IV-C1 and IV-C2, respectively.

A. Dataset

As discussed in Section II, the PV system is simulated in MATLAB Simulink (2021a) and we assume a 6-dimensional waveform data vector is sampled from WMU with a 20kHz sampling rate. We conducted simulations for 9 cases, considering normal conditions along with a variety of attack parameters and irradiation levels. Among them, there are 55,322 normal cases. For cyber attacks 1 through 6, there are 390 cases each. Additionally, there are 17,849 cases pertaining to fault 1 and 7,124 cases associated with fault 2. The normal, fault 1, and fault 2 cases are used to train our FSL model.

B. Model Definition

We used CNN with residual block to train the embedding (embedding layers in Fig. 5). Initially, the model employs a 1-dimensional convolution layer, batch normalization, a ReLU activation function, and max-pooling. Sequentially connected, six residual blocks follow, each comprising two sets of batch normalization, ReLU activation, dropout, and 1-dimensional convolution layers, intertwined with a skip connection that employs max-pooling for shape alignment. A flattening layer subsequently precedes two fully connected layers, the final one matching our number of classes, concluding the model structure. Thus, this structure effectively combines convolutional layers' feature extraction capabilities with residual blocks' ability to leverage multi-layer feature representations, tailoring itself as a potent model for the following classification. Due to page constraints, the detailed structure is illustrated in [15].

C. Experiment Results

In this subsection, we delineate the findings from two case studies conducted to evaluate the performance of our proposed methodology.

- 1) Case Study One: We first validate the attack diagnosis performance of the proposed FSL approach with various shots (i.e. samples in the support set). Table I shows the attack detection F_1 score for the 3-shots, 7-shots, 10-shots, individually. The performance improves as the number of shots in the support set increases, allowing for more candidate comparisons. Highlighting the advantages of FSL, we conducted a comparison against prevalent machine learning and deep learning methods, including Support Vector Machine (SVM), Random Forest (RF), Convolutional Neural Network (CNN) and Long Short-term Memory (LSTM). For the machine learning models, features such as magnitude, frequency, and phase angle were extracted from the data segments. The comparative experiments used six distinct training sizes (20%, 30%, 40%, 50%, 60%, 100%), focusing on both attack detection and classification, as detailed in Table II. The table illustrates that FSL surpasses all other methods when the training dataset is relatively small, requiring only 30 samples for each attack type in the support set. When the benchmarks are trained using the complete dataset (100%), FSL, despite utilizing significantly fewer samples, achieves results that are comparable to most benchmarks. This observation shows the capability of our method to adapt to unseen attacks with only a few examples in the support set.
- 2) Case Study Two: In order to augment the rigor of the experiment, we conducted an analysis of our method's performance using a new case study derived from our previous work [16]. In this work, an online hardware-in-the-loop (HIL) testbed using the OPAL-RT has been built for study on power electronics converters (PECs)-enabled PV farms. The monitoring system runs in real-time while using HIL as an operational solar farm and a National Instruments (NI) data acquisition card as the electric waveform sensor at the point of coupling. Similarly, this work collects data from one PCC node to identify the cyber-attacks and physical faults. The evaluation encompassed various types of cyber-attacks, including denial of service attacks, replay attacks, and fault data injection attacks. Additionally, a normal scenario involving different levels of solar irradiation was considered. The resulting dataset consisted of a total of 91 cases, thereby introducing complexity to the FSL process. In each case, 40 seconds of data samplings of the PV farm are collected, where the sampling rate is 20 kHz. We arbitrarily choose 1 normal case (N1) and 7 attack cases (A1...A7) for building the training set and test set, and another 5 cases (A8...A12) for building the support set and query set. It is noted that the query set is unseen during all training phases.

Table. III displays the attack detection F1 scores, precision and recall for 3-shot, 7-shot, and 10-shot scenarios, further highlighting FSL's effectiveness in performing well with limited samples. To illustrate the advantages of our proposed TSNN implementation for FSL, we also visualize the output embeddings of both the SNN and TSNN models in a 2D space. Figure 6 presents the distinct clustering results obtained from the two implementations. Notably, the clustering results of the TSNN implementation (Figure 6b) surpass those of the SNN

	Detection	Classification	
SVM	0.864/ 0.869/ 0.883/ 0.891/ 0.901/ 0.953	0.535/ 0.568/ 0.587/ 0.601/ 0.642/ 0.903	
RF	0.623/ 0.655/ 0.678/ 0.682/ 0.707/ 0.861	0.525/ 0.533/ 0.549/ 0.571/ 0.618/ 0.827	
. KNN	0.579/ 0.593/ 0.628/ 0.631/ 0.637/ 0.789	0.565/ 0.573/ 0.593/ 0.603/ 0.629/ 0.702	
LSTM	0.826/ 0.844/ 0.879/ 0.885/ 0.901/ 0.958	0.791/ 0.800/ 0.833/ 0.835/ 0.843/ 0.914	
CNN	0.834/ 0.853/ 0.887/ 0.893/ 0.910/ 0.965	0.798/ 0.808/ 0.840/ 0.842/ 0.850/ 0.920	
FSL (30-shots)	0.937	0.899	

TABLE III
DETECTION ON STUDY CASE TWO (F1/ RECALL/ PRECISION)

Support set	3-shots	7-shots	10-shots
attack 8	0.895/ 0.900/ 0.890	0.805/ 0.850/ 0.765	0.824/ 0.880/ 0.775
attack 9	0.955/ 0.960/ 0.950	0.960/ 0.970/ 0.950	0.976/ 0.980/ 0.972
attack 10	0.835/ 0.860/ 0.811	0.945/ 0.950/ 0.940	0.998/ 0.999/ 0.997
attack 11	0.980/ 0.985/ 0.975	0.980/ 0.983/ 0.977	0.988/ 0.992/ 0.984
attack 12	0.960/ 0.975/ 0.945	0.904/ 0.920/ 0.889	0.932/ 0.940/ 0.924

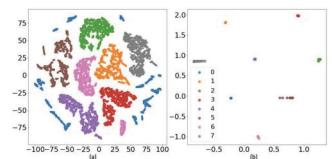


Fig. 6. Comparison of 2D clustering results between SNN implementation (a) and our proposed TSNN implementation (b). The legend is shared between the two subfigures. In the SNN implementation, triplet loss is employed, while the proposed weighted loss function is utilized in the TSNN implementation. implementation (Figure 6a). While the SNN can group samples from the same class together, it struggles to effectively separate them from samples belonging to different classes. Conversely, the TSNN demonstrates improved performance by maintaining both low coupling (distinct separation between classes) and high cohesion (tight clustering of samples within each class). This enhanced clustering capability of the TSNN highlights its efficacy in facilitating accurate discrimination between different classes, thereby enhancing the FSL performance.

V. Conclusions

In this paper, we investigate an under-explored research area in power electronics systems in which only limited labeled and imbalanced data samples are available for cyber-attack diagnosis. We developed an FSL-based approach, which learns the similarity distance among attack samples according to their optimized feature representation. Our proposed framework is capable of 1) dealing with the issue of lacking labeled data for model training, and 2) adapting quickly to detect new attacks without the need to retrain the model. Extensive experimental results demonstrate the effectiveness and efficiency of our proposed approach in detecting cyber-attacks in solar farms and potentially other power electronics systems.

ACKNOWLEDGMENT

This research was partially supported by the U.S. Department of Energy's Solar Energy Technology Office under award

number DE-EE0009026, U. S. National Science Foundation ECCS-EPCN #2102032 and NSF-SATC-2019311, and U. S Department of the Air Force FA8571-20-C-0017.

REFERENCES

- X. Liang, S. A. Wallace, and D. Nguyen, "Rule-based data-driven analytics for wide-area fault detection using synchrophasor data," *IEEE Transactions on Industry Applications*, vol. 53, no. 3, pp. 1789–1798, 2016.
- [2] B. Wang, H. Wang, L. Zhang, D. Zhu, D. Lin, and S. Wan, "A data-driven method to detect and localize the single-phase grounding fault in distribution network based on synchronized phasor measurement," *EURASIP Journal on Wireless Communications and Networking*, vol. 2019, no. 1, p. 195, 2019.
- [3] S. Xingjian, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo, "Convolutional lstm network: A machine learning approach for precipitation nowcasting," in *Advances in neural information processing systems*, 2015, pp. 802–810.
- [4] I. Perry, L. Li, C. Sweet, S.-H. Su, F.-Y. Cheng, S. J. Yang, and A. Okutan, "Differentiating and predicting cyberattack behaviors using lstm," in 2018 IEEE Conference on Dependable and Secure Computing (DSC). IEEE, 2018, pp. 1–8.
- [5] F. Li, Q. Li, J. Zhang, J. Kou, J. Ye, W. Song, and H. A. Mantooth, "Detection and diagnosis of data integrity attacks in solar farms based on multilayer long short-term memory network," *IEEE Transactions on Power Electronics*, vol. 36, no. 3, pp. 2495–2498, 2020.
- [6] L. Zhao, J. Li, Q. Li, and F. Li, "A federated learning framework for detecting false data injection attacks in solar farms," *IEEE Transactions* on Power Electronics, vol. 37, no. 3, pp. 2496–2501, 2021.
- [7] L. Bilge and T. Dumitraş, "Before we knew it: an empirical study of zero-day attacks in the real world," in *Proceedings of the 2012 ACM* conference on Computer and communications security, 2012, pp. 833– 844.
- [8] G. S. Dhillon, P. Chaudhari, A. Ravichandran, and S. Soatto, "A baseline for few-shot image classification," arXiv preprint arXiv:1909.02729, 2019
- [9] E. Schwartz, L. Karlinsky, J. Shtok, S. Harary, M. Marder, A. Kumar, R. Feris, R. Giryes, and A. Bronstein, "Delta-encoder: an effective sample synthesis method for few-shot object recognition," *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [10] Q. Li, F. Li, J. Ye, and W. Song, "Data-driven cyberattack detection for photovoltaic (pv) systems through analyzing micro-pmu data," in *The Twelfth Annual Energy Conversion Congress and Exposition*. IEEE, 2020.
- [11] Q. Li, J. Zhang, J. Zhao, J. Ye, W. Z. Song, and F. Li, "Adaptive hierarchical cyber attack detection and localization in active distribution systems," *IEEE Transactions on Smart Grid*, 2022.
- [12] Y. Wang, Q. Yao, J. T. Kwok, and L. M. Ni, "Generalizing from a few examples: A survey on few-shot learning," *ACM computing surveys (csur)*, vol. 53, no. 3, pp. 1–34, 2020.
 [13] X. Dong and J. Shen, "Triplet loss in siamese network for object
- [13] X. Dong and J. Shen, "Triplet loss in siamese network for object tracking," in *Proceedings of the European conference on computer vision* (ECCV), 2018, pp. 459–474.
- [14] D. Chicco, "Siamese neural networks: An overview," Artificial Neural Networks, pp. 73–94, 2021.
- [15] Q. Li, "Tsnn structure," https://gist.github.com/aaronli-uga/ 958344baeca936ef2efffec51c3a6172, 2024.
- [16] L. Guo, J. Zhang, J. Ye, S. J. Coshatt, and W. Song, "Data-driven cyber-attack detection for pv farms via time-frequency domain features," *IEEE transactions on smart grid*, vol. 13, no. 2, pp. 1582–1597, 2021.