# Learning from the Good Ones: Risk Profiling-Based Defenses Against Evasion Attacks on DNNs

Mohammed Elnawawy
*Department of Electrical and Computer Engineering*
*University of British Columbia*
Vancouver, Canada
mnawawy@ece.ubc.ca

Gargi Mitra
*Department of Electrical and Computer Engineering*
*University of British Columbia*
Vancouver, Canada
gargi@ece.ubc.ca

Shahrear Iqbal
*National Research Council Canada*
Canada
shahrear.iqbal@nrc-cnrc.gc.ca

Karthik Pattabiraman
*Department of Electrical and Computer Engineering*
*University of British Columbia*
Vancouver, Canada
karthikp@ece.ubc.ca

*Abstract*—Safety-critical applications such as healthcare and autonomous vehicles use deep neural networks (DNN) to make predictions and infer decisions. DNNs are susceptible to evasion attacks, where an adversary crafts a malicious data instance to trick the DNN into making wrong decisions at inference time. Existing defenses that protect DNNs against evasion attacks are either static or dynamic. Static defenses are computationally efficient but do not adapt to the evolving threat landscape, while dynamic defenses are adaptable but suffer from an increased computational overhead. To combine the best of both worlds, in this paper, we propose a novel risk profiling framework that uses a risk-aware strategy to selectively train static defenses using victim instances that exhibit the most resilient features and are hence more resilient against an evasion attack. We hypothesize that training existing defenses on instances that are less vulnerable to the attack enhances the adversarial detection rate by reducing false negatives. We evaluate the efficacy of our risk-aware selective training strategy on a blood glucose management system that demonstrates how training static anomaly detectors indiscriminately may result in an increased false negative rate, which could be life-threatening in safety-critical applications. Our experiments show that selective training on the less vulnerable patients achieves a recall increase of up to 27.5% with minimal impact on precision compared to indiscriminate training.

*Index Terms*—risk profiling, evasion attacks, anomaly detectors, selective training, and blood glucose management.

## I. INTRODUCTION

Deep neural networks (DNNs) have gained traction in safety-critical applications such as healthcare [1]–[5] and autonomous vehicles (AVs) [6]–[8]. However, DNNs are highly susceptible to adversarial attacks [9]–[11], especially evasion attacks [12]–[14], which are prevalent since they are relatively easy to execute during deployment [15], [16]. In evasion attacks, a DNN is tricked into misclassifying an adversarial sample at inference time, leading to poor accuracy [17], [18]. For example, an adversary may target DNN models that predict blood glucose values to cause insulin overdose or underdose by manipulating patients' vital signs like previous blood glucose values or administered insulin dosage while ensuring the resulting glucose is within physiological limits to evade detection, leading to catastrophic consequences [19].

Researchers have proposed defense strategies to make DNNs resilient against evasion attacks including adversarial training [20]–[22], training dataset strengthening [23]–[26], model algorithm enhancement [27]–[30], and anomaly detectors [31]–[37]. These defenses are either static or dynamic in nature. Static defenses are easier to implement, demonstrate higher accuracy on benign data, and are more computationally efficient. However, they cannot adapt to different attack strategies or the evolving behavior of victim instances [38]. Dynamic defenses, on the other hand, are more robust to evasion attacks because they adapt to evolving attack and victim behaviors. However, they suffer from degradation of benign data accuracy and high computational overhead at inference time. Therefore, they are not suitable for time-sensitive safety-critical applications [39].

To bridge the gap between static and dynamic defenses, we propose a novel risk-aware selective training strategy that improves the adaptability of static defenses, while retaining their computational efficiency in the presence of an attack. Our risk-aware strategy is powered by a risk profiling framework that selects training instances that show more resilience to the attack. *The key idea is that instances that are less vulnerable to the evasion attack are usually a better representation of a typical distribution of benign data.* As a result, training static defenses to recognize a better distribution of benign data makes it easier for the defense technique to recognize malicious patterns generated by evasion attacks.

In this work, we focus on training anomaly detectors for attack detection in safety-critical applications, e.g., $k$NN, OneClassSVM, and MAD-GAN [31]. The main issue with existing anomaly detectors is that they are often indiscriminately trained on the entire dataset to capture the full spectrum of possible risk scenarios [40]–[42]. However, this strategy has three major problems. *First*, it often yields detectors that are less robust against evasion attacks due to the presence of noisy

data samples which obscure learning meaningful patterns for malicious data detection [43], [44]. *Second*, it degrades the model's generalizability, which is crucial for deploying models in diverse adversarial settings [45]. *Third*, it incurs increased computational cost during training [46]. *We hypothesize that training anomaly detectors using less vulnerable instances can improve malicious data detection by lowering the false negative rate.* We prioritize lower false negatives since higher false negatives in safety-critical applications may lead to deadly consequences whereas higher false positives may lead to denial of service attacks and lack of availability, which are less severe in such systems.

Thus, our goal is to maximize the recall of existing anomaly detectors without causing much degradation to their precision. Towards this goal, we introduce a risk profiling framework that selectively trains existing anomaly detectors on the most resilient instances to help them better differentiate between benign and malicious samples. This boosts the detection rate while overlooking noisy samples that impede the learning process. Our risk profiling framework consists of five steps. *First*, it simulates the evasion attack. *Second*, it quantifies the risk of manipulating data points at every point in time. *Third*, it constructs a time-series risk profile for every victim. *Fourth*, it groups risk profiles depending on their level of vulnerability to the attack. *Fifth*, it uses instances that are less vulnerable to the attack to selectively train the anomaly detectors.

We evaluate the efficacy of our proposed risk profiling framework on a blood glucose management system (BGMS) exposed to evasion attacks against Type-1 diabetes patients. In the context of a BGMS, we define evasion attacks as intentional glucose manipulations designed to deceive DNNs into predicting future glucose levels that result in an altered patient diagnosis. We use the OhioT1DM dataset [47] which includes physiological measurements of 12 Type-1 diabetes patients (six from 2018 and six from 2020). We also use a blood glucose prediction model from prior work [48], which predicts future blood glucose values.

In summary, the contributions of this paper are twofold:

1) A risk profiling framework to quantify the risk of an evasion attack on victims of safety-critical applications and group them into different vulnerability clusters.
2) A strategy to selectively train anomaly detectors on instances with the most resilient features against the attack as identified by the risk profiling framework.

The results of our experiments show that compared to indiscriminate training, selective training guided by our risk profiling framework achieves a recall increase of 27.5% and 16.8% on *k*NN and OneClassSVM, respectively, with little to no impact on precision. Furthermore, when trained on the less vulnerable patients, a MAD-GAN detector maintains a false negative rate of zero with no change to its precision, at a 75% reduction in training set size as opposed to indiscriminately training it on the entire dataset. Therefore, our risk profiling framework helps static anomaly detectors achieve lower false negatives with minimal impact on false positives.
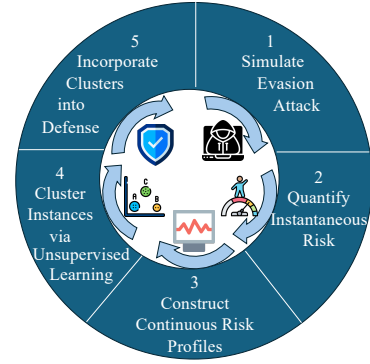


Fig. 1. The five steps of the proposed risk profiling framework.

## II. PROPOSED FRAMEWORK

In this section, we present our risk profiling framework for selective training of existing anomaly detectors to improve their detection capabilities. We rely on anomaly detectors that work in conjunction with the main DNN prediction model. The main DNN model remains unmodified since our proposed risk profiling framework is only used to train the anomaly detectors used to defend against adversarial attack samples. The key idea of the proposed framework is to identify instances that are more resilient due to their natural physiology or driving habits. To do so, the proposed framework categorizes victim instances into clusters of different risk levels depending on their vulnerability to the evasion attack. Once it determines the most resilient instances, the framework uses their past data to selectively train anomaly detectors to recognize the robust features that allow the instances to combat the evasion attack.

Figure 1 shows the proposed risk profiling framework, which consists of five steps. *First*, the framework simulates the evasion attack by generating manipulated inputs to deceive the main DNN model and evaluate its vulnerabilities. *Second*, it quantifies the amount of risk imposed on a victim by calculating the risk metrics at each time stamp to assess the impact of adversarial manipulations on individual data points. *Third*, it constructs a continuous risk profile for each victim to capture their temporal patterns. The risk profile is a time-series representation of all risk values calculated in step 2. *Fourth*, it uses unsupervised machine learning techniques to cluster time-series risk profiles into distinct risk categories, enabling differentiation of vulnerability levels. *Fifth*, it incorporates clustering insights by selectively training anomaly detectors on the less vulnerable instances to learn robust features that improve resilience against evasion attacks.

## III. BLOOD GLUCOSE MANAGEMENT SYSTEM

To evaluate the efficacy of our proposed risk profiling framework in enhancing the performance of existing anomaly detectors, we adopt and extend the case study presented in Elnawawy et al. [49] to simulate evasion attacks. We consider a BGMS (shown in Figure 2) that consists of a continuous glucose monitor (CGM) that measures glucose at regular intervals and transmits it to a smart app running on a mobile
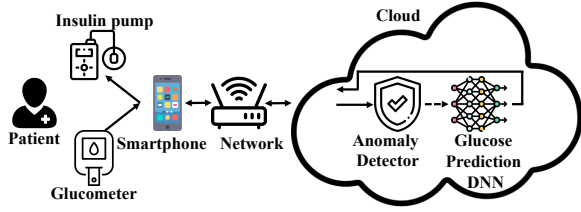
Fig. 2. A BGMS that uses a glucometer, insulin pump, DNN for insulin recommendations, and an anomaly detector to detect adversarial samples.

device via Bluetooth. The app sends the measured glucose to the cloud, where an anomaly detector inspects glucose samples to flag any malicious patterns. If a glucose sample is deemed to be benign by the anomaly detector, it is used by the main DNN model for processing and future glucose predictions. Next, the DNN model sends the predicted future glucose to the mobile app, which calculates the recommended insulin dose and enables the patient to approve it before the insulin pump infuses the corresponding insulin into his/her body.

**Threat model.** The attacker's *goal* is to deceive the BGMS into mistakenly recommending an excessively high insulin dose which could lead the patient into a coma or even death. The attacker's *strategy* is to cause the glucose prediction DNN to predict a high future blood glucose level (hyperglycemia), when in reality the patient has a low (hypoglycemia) or normal blood glucose. To do so, the adversary manipulates the victim's blood glucose levels to values that exceed 125 mg/dL (hyperglycemic in a fasting state) or 180 mg/dL (hyperglycemic two hours postprandial). We assume minimal *capabilities* where the adversary can only manipulate the CGM measurements by compromising the Bluetooth stack via known exploits [50], [51] to intercept and manipulate glucose measurements since many CGM devices use Bluetooth to transmit unencrypted glucose values [52]. Manipulating other features remains beyond the attacker's capabilities. However, we assume the adversary can compromise the smartphone [53] to read these features and ensure the soundness of the generated adversarial samples.

**Target glucose model.** Since the glucose prediction algorithm used by smart apps is often confidential [54], we approximated it using a time-series prediction model developed by Rubin-Falcone et al. [48], which uses a bidirectional long short-term memory (LSTM) architecture. Rubine-Falcone et al. [48] built two types of models: (i) a personalized model for each patient trained on the patient's individual data, and (ii) an aggregate model trained on the data of all patients. We use both types of models to simulate the evasion attack.

**Dataset.** To demonstrate the effect of adversarial glucose values on the target model's predictions, we use the OhioT1DM dataset [47], which was also used by the target model [48] to evaluate its accuracy. The dataset comprises physiological measurements of 12 Type-1 diabetes patients (six from 2018 and six from 2020). For the rest of this paper, we refer to the 2018 and 2020 patients as *Subset A* and *Subset B*, respectively. The main features are the CGM measurements, finger-based measurements, basal insulin, bo-

lus dose, carbohydrate intake, heart rate, sleeping patterns and acceleration, besides other physiological, and self-reported life-event features. The dataset spans eight weeks and consists of ≈10000 samples for training, and 2500 samples for testing, recorded at approximately five-minute intervals per patient.

**Attack algorithm.** As for the evasion attack, we use the universal robustness evaluation toolkit (URET), which is a general-purpose evasion attack framework for manipulating data points at inference time [55]. To ensure that manipulated CGM values respect physiological levels, we constrain them to be between 125 and 499 mg/dL for fasting scenarios, since a hyperglycemic glucose level in a fasting state is greater than 125 mg/dL, and between 180 and 499 mg/dL for postprandial scenarios, since a hyperglycemic glucose level in a postprandial state is greater than 180 mg/dL (499 mg/dL is the highest reported glucose level in the OhioT1DM dataset).

**Anomaly detectors.** To test our framework, we use three anomaly detectors, *k*NN, OneClassSVM, and MAD-GAN [31]. We use *k*NN, for its strength in handling sparse neighborhoods [56], which better represent anomalies in medical data [57]–[59], OneClassSVM for its strength in learning decision boundaries near benign data, making it effective for detecting rare or unusual patterns [60], and MAD-GAN for its strength in capturing multivariate time-series feature dependencies, which is well suited for safety-critical applications like healthcare and AVs [31].

## IV. EVALUATION

In this section, we ask the following research questions:

**RQ1:** Does indiscriminate training of anomaly detectors result in a higher false negative rate? If so, when?

**RQ2:** What is the most suitable selective training strategy to prioritize lower false negatives in anomaly detectors?

To answer the questions, we apply our proposed risk profiling framework to the BGMS discussed in Section III and show how it can be used to enhance the performance of static anomaly detectors using selective training.

**Step 1: Attack Simulation.** In their demonstration of the URET evasion attack on the OhioT1DM dataset using the attack settings presented earlier, Elnawawy et al. [49] show that patients respond differently to the same attack settings as they show different vulnerability levels to the attack. In particular, Elnawawy et al. [49] report attack success rates of mispredicting normal glucose as high glucose reaching up to 100.0% while fasting, and 97.9% postprandial on some patients of *Subset B*, while others show success rates of only 67.4% while fasting, and 44.2% postprandial. This suggests that it is more challenging for URET to attack specific patients who show more resilience to the attack [49]. We extended their experiments to test the URET attack on *Subset A* (Appendix A). The results confirm that different patients of *Subset A* also show different vulnerability levels to the same evasion attack.

**Steps 2 and 3: Risk Quantification.** To quantify the instantaneous risk of an attack at every timestamp, our risk formula considers two factors: (1) magnitude of deviation, and (2) severity/cost of deviation, between the benign and adversarial

TABLE I
SEVERITY COEFFICIENTS FOR DIFFERENT STATE TRANSITIONS

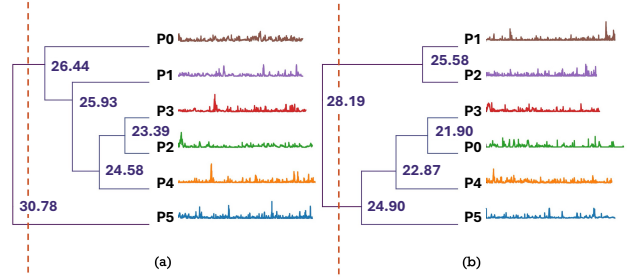| Benign | Adversarial | Severity Coefficient (S) |
|--------|-------------|--------------------------|
| Hypo | Hyper | 64 |
| Normal | Hyper | 32 |
| Hypo | Normal | 16 |
| Hyper | Hypo | 8 |
| Hyper | Normal | 4 |
| Normal | Hypo | 2 |



Fig. 3. The results of hierarchically clustering the risk profiles from (a) *Subset A* and (b) *Subset B* of the OhioT1DM dataset. Based on the distance between the clusters, the dendrograms show that patients in either *Subset* can be clustered into two groups - less and more vulnerable to the attack.

model predictions. The magnitude of deviation is essential for the risk formula since it determines the prediction's state transition. For example, modifying the blood glucose prediction from 90 mg/dL to 210 mg/dL transitions a patient from a state of normal glucose to a state of hyperglycemic glucose. The severity of deviation is important since it weighs state transitions differently depending on the threats they pose to victim instances. For example, transitioning a diabetic patient from hypoglycemic to hyperglycemic glucose is more life-threatening than from normal to hyperglycemic glucose.

In our case study, we calculate the instantaneous risks of manipulating blood glucose values using Equation 1:

$$R_t = S * Z_t, \qquad t \subset \mathbb{N} \qquad (1)$$

where $R_t$ is the instantaneous risk at time unit $t$, $S$ is the severity/cost coefficient of mispredicting a patient's blood glucose level, and $Z_t$ is the difference in magnitude between the benign and adversarial glucose predictions at time unit $t$. $Z_t$ can be calculated using Equation 2:

$$Z_t = (y_t - f(x_t))^2 \qquad (2)$$

where $y_t$ is the benign glucose prediction at time $t$, and $f(x_t)$ is the glucose prediction at time $t$ in the presence of an attack. The difference between $y_t$ and $f(x_t)$ is squared in Equation 2 to weigh big errors more compared to small ones (inspired by the mean squared error) since larger glucose differences could lead to more serious conditions. Next, after the framework calculates instantaneous risk values, it combines them to generate a continuous time-series risk profile for every victim.

Ideally, severity coefficients should be determined by specialists. However, we did not have access to such specialists. Hence, we used exponential coefficients since in healthcare contexts such as BGMS, state transitions (e.g., hypoglycemia to hyperglycemia) are inherently nonlinear in their impact on patient outcomes [61]–[63]. Hence, exponential coefficients capture this nonlinearity by assigning disproportionately higher coefficients to more severe state transitions. Table I shows an example of severity coefficients assigned to different state transitions. For instance, a severity coefficient of 64 is assigned to a diagnosis of hyperglycemia when the actual state of the patient is supposedly hypoglycemic. Hypoglycemia to hyperglycemia misdiagnosis is considered to be the worst case since the system would mistakenly predict an excessively high insulin dose, which could lead to fatal outcomes [64]–[66].

**Step 4: Clustering.** Once the framework generates patients' risk profiles, it uses hierarchical clustering to identify less vulnerable and more vulnerable patients to the attack. In our case study, we chose hierarchical clustering for three reasons [67]. *First*, we do not need to specify the number of clusters in advance since it is difficult to know apriori. Instead, the resulting dendrogram can be pruned at the desired level according to the distances between clusters. *Second*, the dendrogram helps to visually observe patients with similar physiological characteristics at different levels of the hierarchy. *Third*, it is suitable for clinical research since it categorizes mixed populations into more homogeneous groups.

Figure 3 shows the time-series risk profiles for each of the six patients from (a) *Subset A* and (b) *Subset B*. It also shows the resulting dendrograms from hierarchically clustering the 12 patients. Based on the maximum distance between clusters in both cases, we decided to split the patients into two clusters: specifically, patients 0, 1, 2, 3, and 4 from *Subset A* belong to one cluster, and patient 5 belongs to the other cluster. Similarly, patients 0, 3, 4, and 5 from *Subset B* belong to one cluster, and patients 1 and 2 belong to another cluster. By cross-checking the resulting clusters with the misclassification percentages due to the attack reported in Elnawawy et al. [49], on *Subset B* and our extended experiments on *Subset A* (Appendix A), we notice that patient 5 from *Subset A* and patients 1 and 2 from *Subset B* (placed in separate clusters by our risk profiling framework) tend to have the lowest misclassification percentage, meaning that these patients were less vulnerable to the URET attack. On the other hand, the rest of the patients showed a relatively higher misclassification percentage, indicating that they were more vulnerable to the attack. These observations enable us to label the clusters according to patients' misclassification percentages as either less or more vulnerable to the URET attack. The obtained clusters are shown in Table II.

To further analyze the obtained clusters, we plot the ratio of normal to abnormal (i.e., hypoglycemic or hyperglycemic) data points in the original benign trace of the 12 patients in Figure 4. We find that patient 5 from *Subset A* and patient 2 from *Subset B*, who belong to the less vulnerable cluster shown in Table II, show the highest benign normal to abnormal

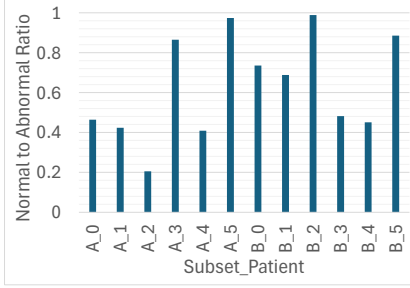| Less Vulnerable | | More Vulnerable | | | | |
|---|---|---|---|---|---|---|
| Subset A | Subset B | Subset A | | | Subset B | |
| p5 | p1 | p0 | p1 | p2 | p0 | p3 |
| | p2 | p3 | p4 | | p4 | p5 |



Fig. 4. Ratio of normal to abnormal data instances in the benign trace of the patients. Less vulnerable patients tend to have higher ratios while more vulnerable patients tend to have lower ratios.



Fig. 5. $k$NN anomaly detection on sample glucose traces of patients 5 and 2 from *Subset A*. Indiscriminately training the detector yields a higher false negative rate on patient 2 (more vulnerable) than patient 5 (less vulnerable).



Fig. 6. The four quadrants of glucose samples: (a) benign normal: normal glucose in absence of attack, (b) benign abnormal: high or low glucose in absence of attack, (c) malicious abnormal: samples intentionally manipulated to fall in the high or low glucose ranges, and (d) malicious normal: samples intentionally manipulated to fall in the normal glucose range.

glucose data points ratio. On the other hand, patient 2 from *Subset A* (more vulnerable cluster) shows the lowest benign normal to abnormal glucose data points ratio.

To demonstrate the issue of indiscriminate training **(RQ1)**, we train a $k$NN anomaly detector using data from all 12 patients of the OhioT1DM dataset. Figure 5 shows sample CGM glucose traces of patients 5 and 2 from *Subset A*. The black and red horizontal lines show the maximum normal glucose values in fasting (125 mg/dL) and postprandial (180 mg/dL) states, respectively. Green dots mark malicious glucose measurements that were successfully flagged by the anomaly detector (i.e., true positives), while red ones mark the missed malicious glucose measurements (i.e., false negatives). The figure shows that indiscriminately training the anomaly detector offers inequitable protection for the two patients since it flagged a higher percentage of adversarial samples from the less vulnerable patient (i.e., patient 5) than the more vulnerable patient (i.e., patient 2). This indicates that the rate of false negatives is much higher for the more vulnerable patient (patient 2) than for the less vulnerable patient (patient 5).

To explain the difference in false negatives caused by indiscriminate training **(RQ1)**, we analyze Figure 4. The more vulnerable patient (A_2) shows a lower ratio of benign normal to abnormal glucose levels, indicating a higher prevalence of abnormal samples in their benign traces (Figure 6). Consequently, when an anomaly detector encounters a malicious abnormal sample, it is more likely to misclassify it as benign because it interprets the abnormality as part of the patient's normal physiological variability rather than a result of an attack, leading to an increased false negative rate. In contrast, the less vulnerable patient (A_5) has a higher ratio of benign normal to abnormal glucose samples, so when the detector sees a malicious abnormal sample there is a higher chance of flagging the sample as malicious (Figure 6). Training on such patients allows the detector to better distinguish between
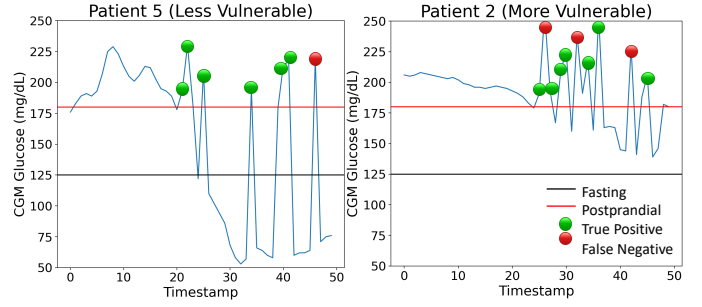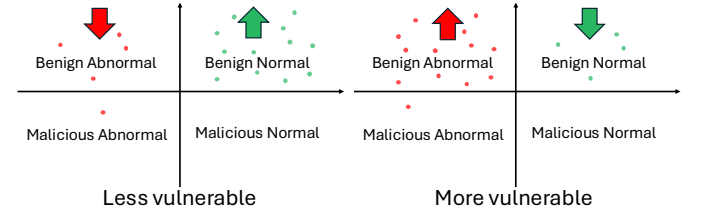
benign and malicious abnormalities, reducing false negatives, albeit at the cost of more false positives (potentially).

**Step 5: Anomaly Detector Enhancement (RQ2).** Training defenses on the less vulnerable patients has three benefits. *First*, it helps learn the most robust features against the attack since the defense focuses on the most resilient, less attack-prone, and more generalizable data features, which helps drive false negatives down. *Second*, it avoids the risk of overfitting to adversarial samples if trained on the more vulnerable patients. This is because when trained on the more vulnerable patients, the model becomes more sensitive to adversarial features, and overfits to specific attack patterns, reducing its generalization ability and driving false positives up. *Third*, the less vulnerable instances better represent a typical distribution of benign data as shown by the higher ratio of benign normal to abnormal glucose samples in Figure 4, providing a more balanced strategy that preserves benign accuracy while detecting attacks. Therefore, we hypothesize that *training defenses on patients who are less vulnerable to the attack enhances model resilience by reducing false negative rate.*

We use four subsets of the OhioT1DM data to train the three anomaly detectors, $k$NN, OneClassSVM, and MAD-GAN [31] (check Appendix B for model parameters). The "Less Vulnerable" and "More Vulnerable" subsets comprise patients shown in Table II. The "Random Samples" subset consists of three patients drawn at random, repeated for 10 different runs, and averaged to reduce random errors and

improve the accuracy of the results. Since the less vulnerable subset is used to test our hypothesis of selectively training anomaly detectors, we randomly sampled three patients in each run of the "Random Samples" experiments to test whether the improvement in less vulnerable training (which included exactly three patients) was purely due to chance. Finally, to show the improvement of selective training on the less vulnerable instances, we indiscriminately train the three defenses on "All Patients" to evaluate the efficacy of our selective training strategy. We consider "All Patients" and "Random Samples" to be our baseline training strategies since they train the anomaly detectors in the absence of insights from our risk profiling framework.

We first consider the recall of the detectors under selective training. Figure 7 shows the results of the recall achieved for each of the training sets. *We observe that training the three defenses on the less vulnerable patients achieves the highest recall among all the subsets, for all three detectors.* In the case of $k$NN and OneClassSVM, training using the less vulnerable patients shows a significant improvement compared to indiscriminately training using the entire dataset, achieving a percentage increase of 27.5% and 16.8% on $k$NN and OneClassSVM, respectively. Moreover, the recall of the less vulnerable training surpasses that of the more vulnerable or randomly sampled patients. In the case of MAD-GAN, training on the less vulnerable patients achieves the same recall as training on the entire dataset (recall of 1), albeit at a 75% reduction of training set size, ensuring better scalability with large, high-dimensional, non-linear, or complex datasets.

We consider the detectors' precisions under selective training. The precision results of $k$NN shown in Figure 8 demonstrate the trade-off between false negatives and false positives since an increase in $k$NN's recall comes at the expense of a 5% reduction in precision with the less vulnerable training. On the other hand, OneClassSVM shows a 7.5% increase in precision when trained using the less vulnerable patients. This may be attributed to $k$NN's sensitivity to the data distribution since glucose data is non-uniformly distributed or has varying densities in different regions; hence $k$NN may label sparse points as anomalies, despite being valid. Conversely, OneClassSVM is less sensitive to density variations since it creates a global model leading to lower false positives. As for MAD-GAN, all training subsets achieved similar precision. Thus, $k$NN and MAD-GAN suffered a small to no loss in precision, while OneClassSVM's precision improved under selective training.

To further investigate the recall-precision trade-off, we calculate their harmonic mean (i.e., F1-score) (Appendix C). We notice that selective training on the less vulnerable patients significantly improves the performance of $k$NN with an F1-score increase of 7.3% compared to indiscriminate training despite the 5% reduction in precision. This indicates that the combined effect of recall and precision captured by their harmonic mean shows an increase in anomaly detection performance, which highlights the efficacy of the proposed framework. On the other hand, OneClassSVM shows an F1-score increase of 10.9% compared to indiscriminate training. The results
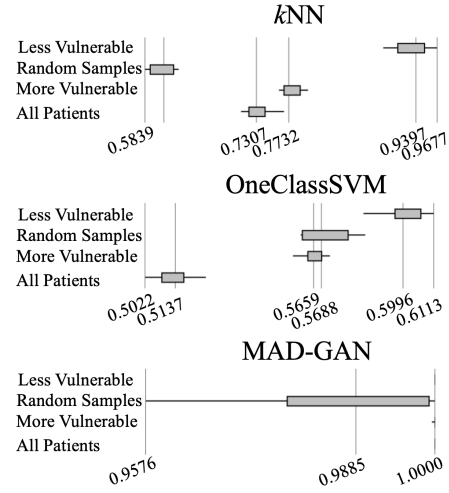


Fig. 7. Recall results using $k$NN, OneClassSVM, and MAD-GAN. Less vulnerable training achieves a recall increase of 27.5% ($k$NN), and 16.8% (OneClassSVM) over indiscriminate training.
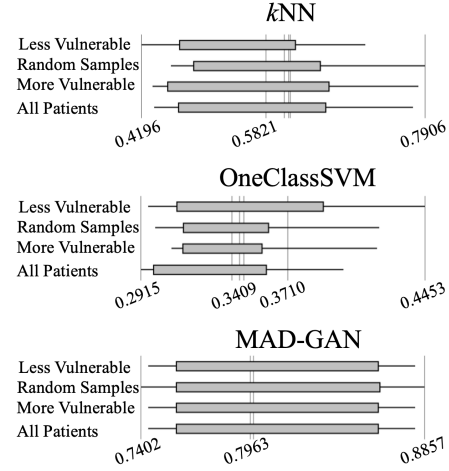


Fig. 8. Precision results using $k$NN, OneClassSVM, and MAD-GAN. Less vulnerable training yields a precision drop of 5% ($k$NN), and an increase of 7.5% (OneClassSVM) over indiscriminate training.

show that despite potential increases in the false positive rate resulting from the recall-precision trade-off, selective training offers an improvement in the combined adversarial detection rate.

## V. LIMITATIONS AND FUTURE WORK

Our framework has four main limitations. *First*, it assumes training and testing data are drawn from the same distribution, which does not consider concept drifts [68]. This leads to a failure to generalize to different data distributions and a failure to adapt to varying environments. For example, a risk profiler trained on senior patients' data may fail with young ones. *Second*, we use offline training to build a static risk profiler, which does not consider potential future dataset shifts. For example, patients move from high-risk to low-risk categories as they recover from medical conditions after the risk profiler has already classified them. *Third*, we used a single case

study and a single attack algorithm to test the efficacy of our proposed framework. More datasets and algorithms are needed for a more thorough evaluation. In the future, we plan to build a risk profiler that uses online learning to consider varying attack environments, different attack algorithms, and potential dataset shifts to design a more adaptive defense. *Fourth*, our choice of severity coefficients is a direct threat to validity since it may impact the correctness of the risk profiles. In the future, we plan to conduct a sensitivity analysis on coefficient choice to further study this problem.

## VI. CONCLUSION

In this paper, we propose a risk profiling framework that bridges the computational overhead gap between static and dynamic defenses against evasion attacks. The proposed framework enhances the adaptability of static defenses to various threat levels using a novel risk-aware selective training strategy that improves adversarial detection rate. The framework generates time-series risk profiles for every victim and clusters them into different risk categories based on their vulnerabilities to evasion attacks. We show that selectively training static anomaly detectors on the less vulnerable victims enhances their detection rates. We evaluated the proposed framework on a Type-1 diabetes case study. Our results show that selective training surpasses indiscriminate training with a reduction in false negatives across three anomaly detectors ($k$NN, OneClassSVM, and MAD-GAN), achieving a recall increase of up to 27.5% with minimal impact on false positives.

## ACKNOWLEDGEMENTS

## REFERENCES

[1] Abdullah Alanazi. Using machine learning for healthcare challenges and opportunities. *Informatics in Medicine Unlocked*, 30:100924, 2022.

[2] Hafsa Habehh and Suril Gohel. Machine learning in healthcare. *Current genomics*, 22(4):291, 2021.

[3] Jenna Wiens and Erica S Shenoy. Machine learning for healthcare: on the verge of a major shift in healthcare epidemiology. *Clinical infectious diseases*, 66(1):149–153, 2018.

[4] Biswajit R Bhowmik, Shrinidhi Anil Varna, Adarsh Kumar, and Rahul Kumar. Deep neural networks in healthcare systems. In *Machine learning and deep learning in efficacy improvement of healthcare systems*, pages 195–226. CRC Press, 2022.

[5] U.S. FDA. AI-Rad Companion (Cardiovascular). URL: https://www.accessdata.fda.gov/scripts/cdrh/cfdocs/cfpmn/pmn.cfm?ID=K183268, Last accessed: Dec 20, 2023.

[6] Mrinal R Bachute and Javed M Subhedar. Autonomous driving architectures: insights of machine learning and deep learning algorithms. *Machine Learning with Applications*, 6:100164, 2021.

[7] Hengrui Chen, Hong Chen, Ruiyu Zhou, Zhizhen Liu, and Xiaoke Sun. Exploring the mechanism of crashes with autonomous vehicles using machine learning. *Mathematical problems in engineering*, 2021(1):5524356, 2021.

[8] Yongqian Xiao, Xinglong Zhang, Xin Xu, Xueqing Liu, and Jiahang Liu. Deep neural networks with koopman operators for modeling and control of autonomous vehicles. *IEEE transactions on intelligent vehicles*, 8(1):135–146, 2022.

[9] Nina Narodytska and Shiva Prasad Kasiviswanathan. Simple blackbox adversarial attacks on deep neural networks. In *CVPR Workshops*, volume 2, 2017.

[10] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Adversarial attacks on deep neural networks for time series classification. In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2019.

[11] Sajjad Amini, Alireza Heshmati, and Shahrokh Ghaemmaghami. Fast adversarial attacks to deep neural networks through gradual sparsification. *Engineering Applications of Artificial Intelligence*, 127:107360, 2024.

[12] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale. *CoRR*, abs/1611.01236, 2016.

[13] Alesia Chernikova, Alina Oprea, Cristina Nita-Rotaru, and BaekGyu Kim. Are self-driving cars secure? evasion attacks against deep neural networks for steering angle prediction. In *2019 IEEE Security and Privacy Workshops (SPW)*, pages 132–137. IEEE, 2019.

[14] J Dinal Herath, Ping Yang, and Guanhua Yan. Real-time evasion attacks against deep learning-based anomaly detection from distributed system logs. In *Proceedings of the Eleventh ACM Conference on Data and Application Security and Privacy*, pages 29–40, 2021.

[15] Gaudenz Boesch. What is adversarial machine learning? attack methods in 2024, June 2024.

[16] Marco Farinetti. *Evasion attacks against machine-learning based behavioral authentication*. PhD thesis, Politecnico di Torino, 2018.

[17] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.

[18] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.

[19] Tamar Levy-Loboda, Eitam Sheetrit, Idit F Liberty, Alon Haim, and Nir Nissim. Personalized insulin dose manipulation attack and its detection using interval-based temporal patterns and machine learning algorithms. *Journal of Biomedical Informatics*, 132:104129, 2022.

[20] Muhammad Shahzad Haroon and Husnain Mansoor Ali. Adversarial training against adversarial attacks for machine learning-based intrusion detection systems. *Computers, Materials & Continua*, 73(2), 2022.

[21] Minh-Hao Van, Wei Du, Xintao Wu, Feng Chen, and Aidong Lu. Defending evasion attacks via adversarially adaptive training. In *2022 IEEE International Conference on Big Data (Big Data)*, pages 1515–1524. IEEE, 2022.

[22] Weimin Zhao, Sanaa Alwidian, and Qusay H Mahmoud. Adversarial training methods for deep learning: A systematic review. *Algorithms*, 15(8):283, 2022.

[23] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L Yuille, and Kaiming He. Feature denoising for improving adversarial robustness. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 501–509, 2019.

[24] Dongqi Han, Zhiliang Wang, Ying Zhong, Wenqi Chen, Jiahai Yang, Shuqiang Lu, Xingang Shi, and Xia Yin. Evaluating and improving adversarial robustness of machine learning-based network intrusion detectors. *IEEE Journal on Selected Areas in Communications*, 39(8):2632–2647, 2021.

[25] Leo Hyun Park, Jaeuk Kim, Myung Gyo Oh, Jaewoo Park, and Taekyoung Kwon. Adversarial feature alignment: Balancing robustness and accuracy in deep learning via adversarial training, 2024.

[26] Jairo Giraldo, David Urbina, CheeYee Tang, and Alvaro A Cardenas. The more the merrier: adding hidden measurements to secure industrial control systems. In *Proceedings of the 7th Symposium on Hot Topics in the Science of Security*, pages 1–10, 2020.

[27] Amin Ghafouri, Yevgeniy Vorobeychik, and Xenofon Koutsoukos. Adversarial regression for detecting attacks in cyber-physical systems. *arXiv preprint arXiv:1804.11022*, 2018.

[28] Adam Goodge, Bryan Hooi, See-Kiong Ng, and Wee Siong Ng. Lunar: Unifying local outlier detection methods via graph neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 6737–6745, 2022.

[29] Xugui Zhou, Maxfield Kouzel, and Homa Alemzadeh. Robustness testing of data and knowledge driven anomaly detection in cyber-physical systems. In *2022 52nd Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN-W)*, pages 44–51, 2022.

[30] Rong Huang and Yuancheng Li. Adversarial attack mitigation strategy for machine learning-based network attack detection model in power system. *IEEE Transactions on Smart Grid*, 14(3):2367–2376, 2022.

[31] Dan Li, Dacheng Chen, Baihong Jin, Lei Shi, Jonathan Goh, and See-Kiong Ng. Mad-gan: Multivariate anomaly detection for time series data with generative adversarial networks. In *International conference on artificial neural networks*, pages 703–716. Springer, 2019.

[32] Rizwan Hamid Randhawa, Nauman Aslam, Mohammad Alauthman, Muhammad Khalid, and Husnain Rafiq. Deep reinforcement learning based evasion generative adversarial network for botnet detection. *Future Generation Computer Systems*, 150:294–302, 2024.

[33] Cao Phan Xuan Qui, Dang Hong Quang, Phan The Duy, Van-Hau Pham, et al. Strengthening ids against evasion attacks with gan-based adversarial samples in sdn-enabled network. In *2021 RIVF International Conference on Computing and Communication Technologies (RIVF)*, pages 1–6. IEEE, 2021.

[34] Rui Shu, Tianpei Xia, Laurie Williams, and Tim Menzies. Omni: Automated ensemble with unexpected models against adversarial evasion attack. *Empirical Software Engineering*, 27:1–32, 2022.

[35] Islam Elgarhy, Mahmoud M Badr, Mohamed MEA Mahmoud, Mostafa M Fouda, Maazen Alsabaan, and Hisham A Kholidy. Clustering and ensemble based approach for securing electricity theft detectors against evasion attacks. *IEEE Access*, 11:112147–112164, 2023.

[36] Usman Ahmed, Jerry Chun-Wei Lin, and Gautam Srivastava. Mitigating adversarial evasion attacks of ransomware using ensemble learning. *Computers and Electrical Engineering*, 100:107903, 2022.

[37] Shunyao Wang, Ryan KL Ko, Guangdong Bai, Naipeng Dong, Taejun Choi, and Yanjun Zhang. Evasion attack and defense on machine learning models in cyber-physical systems: A survey. *IEEE Communications Surveys & Tutorials*, 2023.

[38] Y Wang, T Sun, S Li, X Yuan, W Ni, E Hossain, and HV Poor. Adversarial attacks and defenses in machine learning-powered networks: A contemporary survey. *arXiv preprint arXiv:2303.06302*, 2023.

[39] Francesco Croce, Sven Gowal, Thomas Brunner, Evan Shelhamer, Matthias Hein, and Taylan Cemgil. Evaluating the adversarial robustness of adaptive test-time defenses. In *International Conference on Machine Learning*, pages 4421–4435. PMLR, 2022.

[40] AKM Iqtidar Newaz, Nur Imtiazul Haque, Amit Kumar Sikder, Mohammad Ashiqur Rahman, and A Selcuk Uluagac. Adversarial attacks to machine learning-based smart healthcare systems. In *GLOBECOM 2020-2020 IEEE Global Communications Conference*, pages 1–6. IEEE, 2020.

[41] Xin Li, Deng Pan, and Dongxiao Zhu. Defending against adversarial attacks on medical imaging ai system, classification or detection? In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 1677–1681. IEEE, 2021.

[42] Byunggill Joe, Akshay Mehra, Insik Shin, and Jihun Hamm. Machine learning with electronic health records is vulnerable to backdoor trigger attacks. *arXiv preprint arXiv:2106.07925*, 2021.

[43] Shivani Gupta and Atul Gupta. Dealing with noise problem in machine learning data-sets: A systematic review. *Procedia Computer Science*, 161:466–474, 2019. The Fifth Information Systems International Conference, 23-24 July 2019, Surabaya, Indonesia.

[44] Haojing Shen, Sihong Chen, Ran Wang, and Xizhao Wang. Adversarial learning with cost-sensitive classes. *IEEE Transactions on Cybernetics*, 53(8):4855–4866, 2022.

[45] Mohammed Alawad, Shang Gao, Xiao-Cheng Wu, Eric B Durbin, Linda Coyle, Lynne Penberthy, and Georgia Tourassi. Adversarial training for privacy-preserving deep learning model distribution. In *2019 IEEE International Conference on Big Data (Big Data)*, pages 5705–5710. IEEE, 2019.

[46] Fengxiang He, Shaopeng Fu, Bohan Wang, and Dacheng Tao. Robustness, privacy, and generalization of adversarial training. *arXiv preprint arXiv:2012.13573*, 2020.

[47] Cindy Marling and Razvan Bunescu. The ohiot1dm dataset for blood glucose level prediction: Update 2020. In *CEUR workshop proceedings*, volume 2675, page 71. NIH Public Access, 2020.

[48] Harry Rubin-Falcone, Ian Fox, and Jenna Wiens. Deep Residual Time-Series Forecasting: Application to Blood Glucose Prediction. In *KDH@ECAI*, pages 105–109, 2020.

[49] Mohammed Elnawawy, Mohammadreza Hallajiyan, Gargi Mitra, Shahrear Iqbal, and Karthik Pattabiraman. Systematically assessing the security risks of ai/ml-enabled connected healthcare systems. *arXiv preprint arXiv:2401.17136*, 2024.

[50] Kasper Rasmussen. BLURtooth: Exploiting Cross- Transport Key Derivation in Bluetooth Classic and Bluetooth Low Energy. In *AsiaCCS*, 2022.

[51] CVE. Medtronic MyCareLink Smart Vulnerability. URL: https://www.cve.org/CVERecord?id=CVE-2020-25183, Last accessed: Mar 10, 2025.

[52] Yuchen Niu and Siew-Kei Lam. Securing automated insulin delivery systems: A review of security threats and protectives strategies. *arXiv preprint arXiv:2503.14006*, 2025.

[53] Guillermo Suarez-Tangil, Juan E Tapiador, and Pedro Peris-Lopez. Compartmentation policies for android apps: A combinatorial optimization approach. In *NSS 2015*, pages 63–77, 2015.

[54] U.S. FDA. DreaMed Advisor Pro. URL: https://www.accessdata.fda.gov/scripts/cdrh/cfdocs/cfpmn/denovo.cfm?id=DEN170043, Last accessed: Nov 30, 2023.

[55] Kevin Eykholt, Taesung Lee, Douglas Schales, Jiyong Jang, and Ian Molloy. URET: Universal Robustness Evaluation Toolkit (for Evasion). In *USENIX Security 23*, pages 3817–3833, 2023.

[56] Ming Zhao, Jingchao Chen, and Yang Li. A review of anomaly detection techniques based on nearest neighbor. In *2018 International Conference on Computer Modeling, Simulation and Algorithm (CMSA 2018)*, pages 290–292. Atlantis Press, 2018.

[57] Amir Adler, Michael Elad, Yacov Hel-Or, and Ehud Rivlin. Sparse coding with anomaly detection. *Journal of Signal Processing Systems*, 79:179–188, 2015.

[58] Durgesh Samariya, Jiangang Ma, Sunil Aryal, and Xiaohui Zhao. Detection and explanation of anomalies in healthcare data. *Health Information Science and Systems*, 11(1):20, 2023.

[59] Ruogu Fang, Tsuhan Chen, Dimitris Metaxas, Pina Sanelli, and Shaoting Zhang. Sparsity techniques in medical imaging. *Computerized medical imaging and graphics: the official journal of the Computerized Medical Imaging Society*, 46(Pt 1):1, 2015.

[60] Kun-Lun Li, Hou-Kuan Huang, Sheng-Feng Tian, and Wei Xu. Improving one-class svm for anomaly detection. In *Proceedings of the 2003 international conference on machine learning and cybernetics (IEEE Cat. No. 03EX693)*, volume 5, pages 3077–3081. IEEE, 2003.

[61] Han Li, Quanzhi Lin, Zhiyuan Jiang, and Guoqiang Zhong. Analyzing the non-linear relationship between fasting blood glucose levels and gensini score in patients with stemi. *Frontiers in Cardiovascular Medicine*, 11:1427567, 2024.

[62] Alice Chan, Lutz Heinemann, Stacey M Anderson, Marc D Breton, and Boris P Kovatchev. Nonlinear metabolic effect of insulin across the blood glucose range in patients with type 1 diabetes mellitus. *Journal of diabetes science and technology*, 4(4):873–881, 2010.

[63] Cleveland Clinic. Somogyi effect: What it is, causes, symptoms & treatment. https://my.clevelandclinic.org/health/diseases/11443-somogyi-effect, January 2023. Accessed: March 31, 2025.

[64] Jean-Francois Yale, Breay Paty, Peter A Senior, Diabetes Canada Clinical Practice Guidelines Expert Committee, et al. Hypoglycemia. *Canadian journal of diabetes*, 42:S104–S108, 2018.

[65] Hypoglycemia (low blood sugar).

[66] Mary Jo DiLonardo and Deanna Altomara. Hypoglycemia (low blood sugar), February 2024.

[67] Joshua Noble. What is hierarchical clustering?, August 2024.

[68] Xiaoge Zhang, Felix TS Chan, Chao Yan, and Indranil Bose. Towards risk-aware artificial intelligence and machine learning systems: An overview. *Decision Support Systems*, 159:113800, 2022.
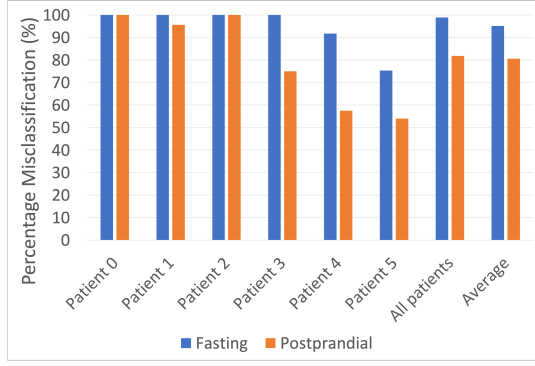
Fig. 9. Percentage of originally normal glucose instances that are misdiagnosed as hyperglycemic. "Patient $i$" shows the results of the personalized model for the $i^{th}$ patient, "All patients" shows the results of the aggregate model trained on the data of all patients, and "Average" shows the average results of the 7 models.
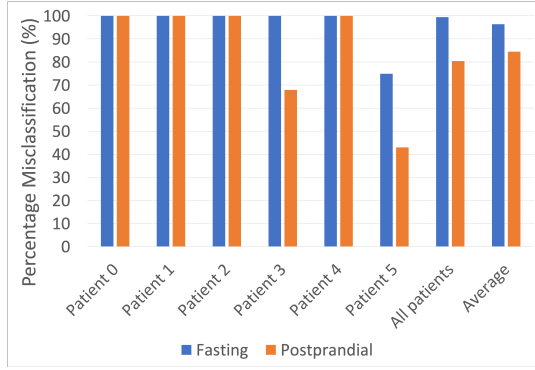


Fig. 10. Percentage of originally hypoglycemic glucose instances that are misdiagnosed as hyperglycemic. "Patient $i$" shows the results of the personalized model for the $i^{th}$ patient, "All patients" shows the results of the aggregate model trained on the data of all patients, and "Average" shows the average results of the 7 models.

In this appendix we elaborate on the anomaly detectors used to test our risk profiling framework.

*k*NN. We use the KNeighborsClassifier implementation of the scikit-learn Python library with the following model parameters:

- Number of neighbors = 7
- Weights = uniform
- Algorithm = auto
- Leaf size = 30
- p = 2
- Metric = minkowski
- Metric params = None

**OneClassSVM.** We use the OneClassSVM implementation of the scikit-learn Python library with the following model parameters:

- Kernel = sigmoid
- Degree = 3
- Gamma = auto
- Coef0 = 10
- Tol = 0.001
- Nu = 0.5
- Shrinking = True
- Cache size = 200
- Max iter = -1

**MAD-GAN [31].** MAD-GAN is an unsupervised anomaly detection technique for multivariate time-series data. It uses a generative adversarial network (GAN) with long short-term memory recurrent neural networks (LSTM-RNN) as the generator and discriminator. MAD-GAN captures temporal correlations and latent interactions among features to detect anomalies using a novel anomaly score called discrimination and reconstruction anomaly score (DR-Score). We use the following model parameters in our adoption of MAD-GAN:

- Number of epochs = 100
- Number of signals = 4
- Number of generated features = 4
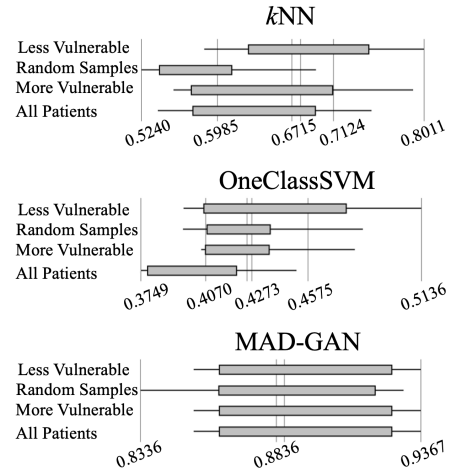- Sequence length = 12
- Sequence step = 1

Fig. 11. F1-score results using *k*NN, OneClassSVM, and MAD-GAN. Less vulnerable training achieves an F1-score increase of 7.3% (*k*NN), and 10.9% (OneClassSVM) over indiscriminate training.

In our experiments, we trained anomaly detectors on the less vulnerable patients and independently tested the entire set of patients and then averaged the results to aggregate them in box plots. That means that we conducted experiments where the test set consisted of only the more vulnerable patients that were not seen during the training stage. The resulting anomaly detection rates were similar to those obtained by testing on the less vulnerable patients. This demonstrates our

framework's resilience to overfitting due to training only on the less vulnerable patients when tested on the OhioT1DM dataset. Nevertheless, we acknowledge that more rigorous testing on different datasets and attack algorithms is needed before confidently claiming the framework's resilience to overfitting and its generalizability to other domains. For this reason, we plan to extend our experiments to other healthcare datasets, the domain of autonomous vehicles (AVs), and other attack algorithms in our next publications. To further validate our work, we believe that our proposed framework could benefit from a mathematical model formulation to capture its full dynamics.

The proposed framework addresses concept drift through an iterative process that regularly reassesses patient risk profiles and continuously updates them as new data become available to strengthen defenses against the evolving threat landscape. As patient conditions evolve, so do their risk levels: those showing increased resilience against adversarial attack are incorporated into the retraining process, while those becoming more vulnerable are excluded from the occasional retraining. This continuous refinement ensures that the defense adapts over time to maintain robustness without sacrificing accuracy.