# `ADAPT`: A Pseudo-labeling Approach to Combat Concept Drift in Malware Detection

1st Md Tanvirul Alam
*Rochester Institute of Technology*
Rochester, NY, USA
ma8235@rit.edu

2nd Aritran Piplai
*University of Texas El Paso*
El Paso, TX, USA
apiplai@utep.edu

3rd Nidhi Rastogi
*Rochester Institute of Technology*
Rochester, NY, USA
nxrvse@rit.edu

*Abstract*—Machine learning models are commonly used for malware classification; however, they suffer from performance degradation over time due to concept drift. Adapting these models to changing data distributions requires frequent updates, which rely on costly ground truth annotations. While active learning can reduce the annotation burden, leveraging unlabeled data through semi-supervised learning remains a relatively underexplored approach in the context of malware detection. In this research, we introduce `ADAPT`, a novel pseudo-labeling semi-supervised algorithm for addressing concept drift. Our model-agnostic method can be applied to various machine learning models, including neural networks and tree-based algorithms. We conduct extensive experiments on five diverse malware detection datasets spanning Android, Windows, and PDF domains. The results demonstrate that our method consistently outperforms baseline models and competitive benchmarks. This work paves the way for more effective adaptation of machine learning models to concept drift in malware detection.

*Index Terms*—malware detection, concept drift, semi-supervised learning, pseudo-labeling

## I. INTRODUCTION

The ever-evolving nature of malware poses ongoing challenges to cybersecurity. While machine learning offers substantial improvements over traditional signature- and rule-based detection [7], [12], [39], [54], [57], its effectiveness depends on the assumption that data distributions remain stable [18]. In reality, continuous changes in malware and underlying platforms—such as new features, APIs, and programming practices [41]—frequently violate this assumption, causing *concept drift* and leading to significant model degradation unless adaptation measures are taken [47], [55].

Dealing with concept drift in malware detection is a multifaceted challenge. While some approaches focus on building more robust feature spaces to resist drift [54], [83], true immunity remains elusive given malware's dynamic landscape [14]. Periodic model retraining offers another solution [50], [55], but requires timely drift detection and access to newly labeled samples [14], [40], [80]. Active learning helps reduce annotation effort [24], [51], yet obtaining high-quality malware labels is costly due to the sheer volume of emerging malware variants [20], [49].

Self-training offers a way to adapt to concept drift without requiring new annotations beyond the initial training data [6]. Pseudo-labeling is a learning paradigm where the model generates labels for unlabeled data, treating its predictions as ground truth to train itself [44], [62] further. A notable application of this approach for malware detection is DroidEvolver [77], which employs pseudo-labeling for Android malware detection by using an ensemble of online learning algorithms to eliminate labeling costs after the initial training phase. However, a later study [41] revealed that when experimental biases, such as unrealistic benign-to-malware ratios in the dataset, are addressed, DroidEvolver's performance deteriorates significantly due to confirmation bias [11] and a catastrophic self-poisoning effect, leading to a sudden and severe performance drop. Although the work in [41] proposed an enhanced version of DroidEvolver to mitigate these issues, performance still deteriorates rapidly after a short period when solely relying on pseudo-labels. Thus, designing a practical pseudo-labeling-based approach for malware detection under concept drift remains a significant challenge.

In this work, we propose `ADAPT`: **A**daptive **D**rift-aware **A**lgorithm using **P**seudo-labeling for malware detec**T**ion, a semi-supervised approach tailored for concept drift in malware detection. Unlike general-purpose methods, `ADAPT` accounts for the asymmetric nature of drift, which more frequently impacts malware than benign samples [14], [27], by introducing a drift-aware pseudo-labeling strategy that selectively propagates labels based on class-specific drift behavior. Data augmentation, though standard in image classification [28], remains uncommon in malware detection, but has shown promise in noisy label settings [76]. We extend these augmentation strategies to the concept drift scenario, periodically retraining the model to adapt to changing distributions. To further enhance robustness, we incorporate mixup regularization [70], [81], which improves confidence calibration and helps mitigate the self-poisoning effect commonly seen in pseudo-labeling. This ensures that the model's predicted probabilities better reflect the actual likelihood of correctness, making pseudo-labeling more effective in practice.

We evaluate `ADAPT` on five diverse malware detection datasets—two Android, two Windows, and one PDF—and show that it consistently outperforms competitive semi-supervised learning methods, achieving state-of-the-art results across all tasks. We also demonstrate its effectiveness in an active learning setting, where `ADAPT` surpasses prior state-of-the-art performance on Android malware detection under the same annotation budget.

Our key contributions are:

1) We introduce ADAPT[1], a semi-supervised learning algorithm tailored for concept drift in malware detection. It mitigates catastrophic self-poisoning through drift-aware pseudo-labeling.

2) Our method is model-agnostic and shows improvements across Random Forest, XGBoost, and neural networks. We validate its effectiveness on five real-world malware datasets exhibiting varying degrees of concept drift.

3) We extend ADAPT to an active learning setting, achieving state-of-the-art results on Android malware detection. We also apply it to the multiclass problem of malware family classification and demonstrate its effectiveness under concept drift.

## II. BACKGROUND AND RELATED WORK

### A. Malware Detection using Machine Learning

Machine learning has been widely applied to malware detection [7], [12], [25], [46], [54], [83], both in academic research and industrial settings. These models are trained using features extracted through static analysis, dynamic analysis, a combination of both, or memory-based analysis [64]. Static features are obtained by analyzing the binary or source code of a file without executing it, making them a popular choice for malware detection [7], [12], [54]. While machine learning detectors built upon static features often achieve high performance when the training and test datasets follow the same distribution, their effectiveness deteriorates in real-world scenarios due to distributional changes caused by code obfuscation, concept drift, and adversarial examples [31]. In particular, concept drift, also known as dataset shift, is the primary focus of this study.

### B. Concept Drift

Dataset shift refers to the phenomenon in machine learning where the joint distribution of inputs and outputs differs between the training and test stages [56]. Dataset shifts can be broadly categorized into three types: covariate shift, label shift, and concept shift. Covariate shift occurs when the distribution of the features, $p(x)$, changes. Label shift, also known as prior probability shift, refers to a change in the distribution of the labels, $p(y)$. Concept shift refers to a change in the conditional probability distribution, $p(y|x)$, meaning the relationship between the features and the target labels changes, effectively altering the class definitions. The individual effects of these shifts can be challenging to disentangle from a finite set of samples [14], [56]. Consequently, in machine learning for security, it is common practice to refer to all these shifts under the broader term *concept drift* [14], [40], [41], which is the terminology we adopt in this work.

Concept drift in malware detection arises from both benign and malicious changes. Updates to APIs and the introduction of new functionalities in benign applications contribute to changes in the input feature space, often resulting in what

---

[1] https://github.com/aiforsec/ADAPT

---

is referred to as *virtual drift* [63]. In contrast, the primary driver of concept drift is the continually evolving behavior of malware [26], [83], which leads to *real concept drift*, a shift in the relationship between features and labels [63]. Adversaries regularly adapt through obfuscation, packing, and novel attack strategies to evade detection [1], [8], [68]. These dynamics alter the distribution of malicious behavior between training and testing data [14], degrading classifier performance and reducing detection rates over time [40], [55].

### C. Concept Drift Adaptation

While developing feature spaces that are inherently more resilient to concept drift is an effective approach [54], [83], the dynamic nature of malware makes it challenging to design such feature spaces that remain effective over long periods [14]. As a result, machine learning models trained on these features need to be periodically retrained to stay functional.

Approaches such as Transcend [40] and Transcendent [14] address concept drift by employing classification with rejection, using conformal evaluators to flag uncertain samples for expert review. However, reliance on manual labeling inherently limits the frequency of model retraining due to cost and resource constraints [49], [77]. Active learning has been proposed to alleviate the annotation burden by selectively labeling only the most uncertain samples [24], [79], with recent studies demonstrating up to an eightfold reduction in labeling costs [24].

Despite these advances, practical challenges such as label delays can significantly increase user exposure to undetected malware [19]. As a result, there is a clear need for interim solutions that can adapt to drift in the absence of timely ground-truth labels. Pseudo-labeling provides one such solution, enabling models to update themselves using confident predictions in place of manual annotations [19], [41], [77].

### D. Pseudo-Labeling for Drift Adaptation

Pseudo-labeling is a technique that leverages model predictions to assign labels to unlabeled data, assuming that predictions with high confidence are likely to be accurate. These pseudo-labels are then incorporated into the training process, enabling the model to learn from labeled and unlabeled samples. Pseudo-labeling is widely used in semi-supervised learning setups [17], [65]. However, in the context of malware detection, its application typically does not account for adapting to concept drift [48], [52], [61], [76].

DroidEvolver [77] applies pseudo-labeling for Android malware detection under concept drift using an ensemble of linear models and online learning to avoid manual labeling. However, it was later shown to underperform when accounting for dataset biases [41], prompting several improvements. More recent work includes MORPH [3], which uses neural networks but suffers from high false positive rates, and Insomnia [9], which applies co-training for drift adaptation in network intrusion detection.
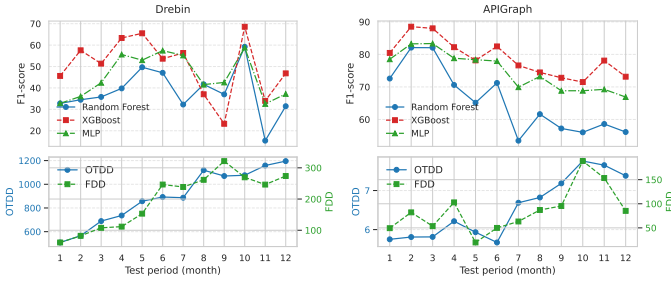
Fig. 1: Covariate shift $p(x)$ and concept shift $p(y|x)$ measured by Optimal Transport Dataset Distance (OTDD) and Fréchet Dataset Distance (FDD) across the first year of testing on two Android datasets (Drebin and APIGraph). The corresponding monthly F1-scores from three models (Random Forest, XG-Boost, and MLP) are shown to illustrate the relationship between distributional shifts and malware detection performance.
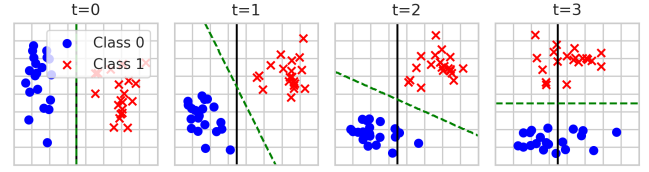


Fig. 2: Adaptation to distribution shift using pseudo-labels. Solid line: original classifier; dotted lines: adapted classifier trained with pseudo-labels.

Despite its potential, pseudo-labeling under concept drift poses challenges: as distributions shift, model predictions become less reliable, leading to incorrect pseudo-labels and self-reinforcing errors—known as self-poisoning [11], [41]. In this work, we propose a pseudo-labeling approach specifically designed to mitigate these effects.

## III. RESEARCH MOTIVATION

*a) Nature of Concept Drift in Malware Detection:* Machine learning models used for malware detection experience degraded performance over time due to concept drift. Although precise measurement of dataset shifts from finite samples is challenging, they can be approximated using metrics such as the Optimal Transport Dataset Distance (OTDD) [5] for covariate shift and the Fréchet Dataset Distance (FDD) [38] for concept shift [33].

OTDD [5] quantifies covariate shift by calculating the minimal cost to transform one data distribution into another using optimal transport theory. Specifically, we approximate this distance with a Gaussian approximation computed directly in the feature space between the training data and each subsequent month's testing data. In contrast, Fréchet Dataset Distance (FDD) [38], initially introduced in the context of generative modeling, captures concept shift by evaluating the similarity of intermediate feature representations extracted from a trained reference classifier. Here, we employ a Multilayer Perceptron (MLP) model trained with optimal hyperparameters (as detailed in Section VI-C) and compare penultimate layer features from each test month's samples against those of the training set.

Applying these metrics to two prominent Android malware datasets, Drebin [12] and APIGraph [83] (Figure 1), we observe substantial and progressively increasing covariate and concept shifts. These distributional shifts are notably more pronounced in Drebin, which utilizes a larger, less semantically compact feature space compared to APIGraph's more concise API-based representation. Although both datasets experience degradation in malware detection accuracy over time,

this deterioration more consistently correlates with increased OTDD and FDD distances in APIGraph, highlighting complex interactions between dataset shifts and classifier performance. These insights underscore the critical importance of continuously adapting models to mitigate the negative impacts of concept drift.

*b) Challenges of Label Delays in Practice:* Effective adaptation to concept drift fundamentally relies on the availability of timely ground-truth labels for evaluating model predictions. However, in practice, these labels often arrive after significant delays due to resource-intensive processes. Sandbox analyses, essential for automated malware evaluation, face scalability constraints relative to the volume of daily submissions, resulting in prolonged analysis queues [30]. Human analysts are even scarcer, causing manual labeling to become a severe bottleneck with substantial waiting periods [15]. Consequently, these delays open extended vulnerability windows, significantly increasing user exposure to emerging threats and degrading detection performance [19].

*c) Pseudo-labeling for Effective Drift Adaptation:* To mitigate these practical challenges associated with label delays, pseudo-labeling emerges as a promising strategy. Pseudo-labeling is a self-training technique in which model-generated predictions serve as provisional labels for unlabeled data, enabling continuous model adaptation without immediate reliance on manual annotations [43]. This process is illustrated in Figure 2 using a 2D toy dataset with a logistic regression classifier. In this example, both $p(x)$ and $p(y|x)$ gradually change from $t = 0$ to $t = 3$, with the dataset undergoing continuous rotation. As a result, the classifier's initial decision boundary (solid line) misclassifies most samples by $t = 3$.

However, by using the model's own predictions as pseudo-labels and retraining on these pseudo-labeled samples, the model can adapt to the changing data distribution. The updated classifier (dashed line) in Figure 2 successfully separates the dataset at the intermediate stages and in the final step. While this example demonstrates the effectiveness of self-training in a simple, idealized scenario, real-world malware datasets are significantly more complex. A naive application of pseudo-labeling in these cases can lead to performance degradation [41].

Nevertheless, malware datasets typically undergo gradual distributional shifts, as illustrated in Figure 1, where both covariate and concept shifts increase over time. Based on

this observation, we hypothesize that an effective self-training approach can continuously adapt the model to account for these evolving shifts. Notably, prior work by Kumar et al. [43] establishes formal error bounds of self-training under the assumption of gradual distributional shifts. Our method builds on this foundation by extending their framework, and we provide a rigorous theoretical analysis in Appendix H.

## IV. PROPOSED METHOD

We introduce `ADAPT`, a pseudo-labeling-based self-training algorithm for handling concept drift in malware detection. `ADAPT` is compatible with various learning algorithms, including tree-based methods and neural networks. We evaluate three baseline models in our experiments: Random Forest (RF), XG-Boost, and Multilayer Perceptron (MLP). These models were selected because different algorithms may perform better on different datasets, and simpler models like Random Forest can be preferable due to their efficiency in specific applications [2], [29], [40], [66], [79].

We assume that we initially have access to an annotated, labeled dataset. Specifically, we are given a set of $N$ annotated samples with features $\{x_1, x_2, \ldots, x_N\}$ and their corresponding binary labels $\{y_1, y_2, \ldots, y_N\}$, where 0 represents a benign sample and 1 represents malware. We first train a supervised machine learning model, denoted as $M_0$, using this labeled data. We later extend our analysis to the multiclass setting, which is discussed in Section IX.

After the initial training, the model is periodically updated each subsequent month using unlabeled test data. This experimental setting aligns with prior work on concept drift adaptation, such as the approach by Chen et al. [24] for continuous learning. We employ our proposed algorithm during each monthly retraining period to leverage the unlabeled data, improving the model's ability to adapt to concept drift.

`ADAPT` comprises four key steps during the semi-supervised update phase:

1) **Adaptive Drift-Aware Pseudo-Labeling:** In this step, we select model predictions that are eligible for inclusion in the self-training process. We propose a novel adaptive pseudo-label selection method to improve the label quality of the selected samples. Details of the pseudo-labeling method are provided in Section IV-A.

2) **Data Augmentation with Label Consistency:** To increase the diversity of the training data, we apply a data augmentation phase. Augmentation has proven to be an effective strategy in self-training and semi-supervised learning, especially in scenarios like malware classification where labels can be noisy. The specifics of our augmentation technique are discussed in Section IV-B.

3) **Confidence Calibration with Mixup:** High-quality pseudo-labels are crucial for successful self-training. Poor pseudo-labels can lead to issues such as self-poisoning, where the model reinforces its incorrect predictions, and confirmation bias. Since we use confidence thresholding to filter out low-confidence predictions, it is essential to calibrate the model's confidence levels—i.e., to ensure

---

**Algorithm 1** `ADAPT`: Adaptive Drift-aware Algorithm using Pseudo-labeling for malware detection

---

1: **Input:** Labeled dataset $\mathcal{D}_l$, unlabeled dataset $\mathcal{D}_u$
2: **Output:** Updated model $M$
3: Train initial model $M_0$ on the labeled dataset $\mathcal{D}_l$
4: **for** each test month from $i = 1$ to $T$ **do**
5:     **Step 1: Pseudo-Label Selection**
6:     Select pseudo-labeled samples: $\mathcal{D}_p = \{(x_i, \hat{y}_i)\}$ from the unlabeled dataset $\mathcal{D}_u$ using the model $M_{i-1}$ [see §IV-A].
7:     Merge the labeled and pseudo-labeled data: $\mathcal{D}_m = \mathcal{D}_l \cup \mathcal{D}_p$
8:     **Step 2: Data Augmentation**
9:     Apply data augmentation to $\mathcal{D}_m$ to generate the augmented dataset $\mathcal{D}_{\text{aug}}$ [§IV-B].
10:     **Step 3: Mixup Regularization**
11:     Apply mixup regularization on $\mathcal{D}_m$ to produce the mixup-augmented dataset $\mathcal{D}_{\text{mixup}}$ [§IV-C].
12:     **Step 4: Model Retraining**
13:     Combine the merged, augmented, and mixup datasets: $\mathcal{D}_{\text{combined}} = \mathcal{D}_m \cup \mathcal{D}_{\text{aug}} \cup \mathcal{D}_{\text{mixup}}$
14:     Retrain the model $M_i$ for the current month on the combined dataset $\mathcal{D}_{\text{combined}}$
15: **end for**

---

that predictions with high confidence are indeed more likely to be correct. To achieve this, we integrate the mixup [81] regularization technique, which has been shown to improve model calibration in prior work [70]. Further details on the mixup method are provided in Section IV-C.

4) **Model Retraining:** In the final step, we combine the original labeled data, the pseudo-labeled data, and the datasets generated through augmentation and mixup regularization. The model is then updated on this combined dataset. For traditional models like random forests and XGBoost, we retrain the model from scratch, while for neural networks, we fine-tune the model, following the guidelines in [24].

Algorithm 1 provides an overview of the `ADAPT` algorithm.

### A. Adaptive Drift-Aware Pseudo-Labeling

In traditional pseudo-label selection, a fixed threshold $\tau$ is used to filter out low-confidence predictions [44]. Given a model $M$ and its predicted probability for a sample $\mathbf{x}_j$, the pseudo-label $\tilde{y}_j$ is defined as:

$$\tilde{y}_j = \begin{cases} \arg\max_k \hat{y}_j^{(k)}, & \text{if } \max_k \hat{y}_j^{(k)} \geq \tau, \\ \text{unlabeled}, & \text{otherwise.} \end{cases}$$

Here, $\hat{y}_j^{(k)}$ denotes the predicted probability of the $k$-th class for the unlabeled sample $\mathbf{x}_j$. This formulation assumes that the model outputs probabilities for two classes, benign and malware, which sum to 1. If the model instead outputs only the probability $p_1$ for the positive (malware) class—such as in
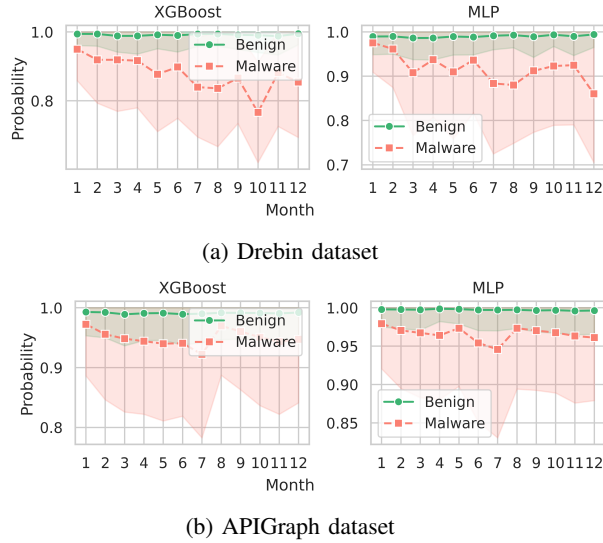
(a) Drebin dataset



(b) APIGraph dataset

Fig. 3: Mean and standard deviation of predicted probabilities for benign and malware classes using XGBoost and MLP. Lower confidence for malware highlights that drift mainly affects this class.

the case of a neural network with a single output neuron and sigmoid activation—then the probability for the benign class is $p_0 = 1 - p_1$. In either case, if the predicted probability for the assigned class exceeds the predefined threshold $\tau$, the sample is pseudo-labeled; otherwise, it remains unlabeled.

Using a fixed thresholding approach for malware detection has limitations due to the divergent drift patterns in the benign and malware classes, as shown in Figure 3. In this figure, we plot the average prediction probability for the benign and malware classes using the XGBoost and MLP models on the Android datasets. The malware class experiences more severe drift, as indicated by its consistently lower average prediction probability than the benign class.

To address this, we propose a class-dependent adaptive thresholding strategy for pseudo-label selection, which dynamically adjusts to the evolving data distribution. While dynamic thresholding has been applied in other domains [74], these approaches do not consider the distinct nature of drift across different classes. In our method, we define separate thresholds for the two classes: $\tau_m \in [0.5, 1]$ for malware and $\tau_b \in [0.5, 1]$ for benign samples. The adaptive thresholding mechanism adjusts these thresholds based on the model's average prediction probabilities for each class over the unlabeled data. Specifically, we compute the mean probability for each class as follows:

$$\mu_m = \frac{1}{|\mathcal{D}_u^m|} \sum_{x_j \in \mathcal{D}_u^m} P(y = 1 \mid x_j) \tag{1}$$

$$\mu_b = \frac{1}{|\mathcal{D}_u^b|} \sum_{x_j \in \mathcal{D}_u^b} P(y = 0 \mid x_j) \tag{2}$$

where $\mathcal{D}_u^m$ and $\mathcal{D}_u^b$ represent the subsets of the unlabeled dataset $\mathcal{D}_u$ that are predicted as malware and benign by the model $M$, respectively. The quantities $\mu_m$ and $\mu_b$ are the mean probabilities for the malware and benign classes. While in our experimental setup, these values are computed using all the samples from the current month's data, the formulation can also be applied in a batch or online learning setting, where the quantities are computed based on batch statistics of the model's predicted pseudo-labels.

We then combine these mean probabilities with the prior thresholds $\tau_m$ and $\tau_b$ using a parameter $\lambda \in [0, 1]$ to compute the updated thresholds, which accounts for the gradual shift in the data:

$$\tau_m^{\text{updated}} = \lambda \cdot \mu_m + (1 - \lambda) \cdot \tau_m \tag{3}$$

$$\tau_b^{\text{updated}} = \lambda \cdot \mu_b + (1 - \lambda) \cdot \tau_b \tag{4}$$

After updating the thresholds, pseudo-labeled data is selected using these new adaptive thresholds. Specifically, if the model's predicted probability for a sample exceeds the updated threshold for malware ($\tau_m^{\text{updated}}$), the sample is assigned a label of 1 (malware). If the probability exceeds the updated threshold for benign ($\tau_b^{\text{updated}}$), it is labeled as 0 (benign). Otherwise, the sample remains unlabeled. Formally:

$$\hat{y}_j = \begin{cases} 1, & \text{if } P_M(y = 1 \mid x_j) > \tau_m^{\text{updated}} \\ 0, & \text{if } P_M(y = 0 \mid x_j) > \tau_b^{\text{updated}} \\ \text{unlabeled}, & \text{otherwise} \end{cases}$$

The values of the threshold parameters $\tau_m$ and $\tau_b$, as well as the adaptation parameter $\lambda$, are jointly optimized during hyperparameter tuning, along with other model parameters, as discussed in Section VI-C.

### B. Data Augmentation with Label Consistency

Data augmentation is a commonly used technique in semi-supervised learning, particularly in contrastive learning, where the model is encouraged to produce consistent predictions for a sample and its perturbed counterpart [17], [65]. While data augmentation methods are prevalent in domains like computer vision, this is not the case for malware detection. In image classification, various label-preserving transformations, such as rotation, translation, scaling, and flipping, can be easily applied [28]. However, designing equivalent transformations for malware datasets is significantly more challenging. Malware datasets typically consist of hand-crafted features, making it difficult to determine whether random transformations in the feature space preserve both the functionality and maliciousness of the application.

To address this challenge, Wu et al. [76] proposed a data augmentation strategy for Windows malware classification in which features are randomly masked and replaced with features from other samples in the dataset. Given a sample $\mathbf{x} \in \mathbb{R}^d$, a mask vector $\mathbf{m} = [m_1, \ldots, m_d]^\top \in \mathbb{R}^d$ is first generated, where each $m_j$ is independently sampled from
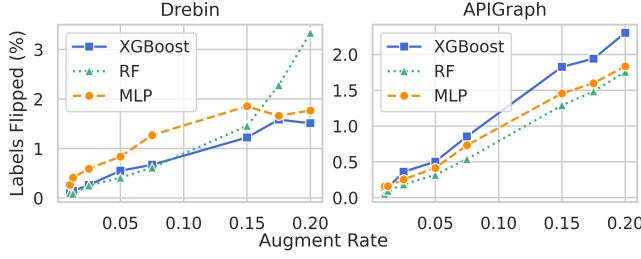
Fig. 4: Percentage of Drebin and APIGraph samples where augmentation flips a correct prediction (i.e., benign to malicious or vice versa).

a Bernoulli distribution with probability $p_a$. The augmented sample $\tilde{\mathbf{x}}$ is then obtained using the following equation:

$$\tilde{\mathbf{x}} = (1 - \mathbf{m}) \odot \hat{\mathbf{x}} + \mathbf{x} \odot \mathbf{m} \tag{5}$$

where $\odot$ denotes element-wise multiplication, and $\hat{\mathbf{x}}$ is a new sample where each feature $\hat{x}_i$ is sampled from the empirical marginal distribution of the $i$-th feature. This empirical marginal distribution is represented by the uniform distribution over the values that the feature takes across the entire training dataset.

For example, consider a dataset with three features and four samples:

$$\{[0, 1, 0], [1, 1, 0], [0, 0, 1], [0, 1, 1]\}.$$

The value of the second feature of $\hat{\mathbf{x}}$ would be randomly sampled from the set $\{1, 1, 0, 1\}$, which are the observed values of the second feature across all samples.

While the data augmentation strategy proposed by Wu et al. [76] was designed for semi-supervised learning in a static setting where the model is trained only once, our scenario involves multiple model updates with additional pseudo-labeled data. Data augmentation can introduce samples whose labels may change due to concept drift in this setting. Such label flipping can degrade model performance in subsequent stages. We illustrate this effect in Figure 4 using the Drebin and APIGraph datasets across different models. Specifically, we select samples from the validation set where the model's original prediction is correct and apply augmentation with varying masking probabilities. We then recompute the model's predictions on the augmented data. Since augmentation should not alter the ground truth, if the model's prediction changes and becomes incorrect, we classify this as a case of label flipping. As the strength of the augmentation increases, the proportion of samples experiencing label flipping also rises. This suggests that overly strong augmentations can cause samples to cross the model's decision boundary, introducing additional noise into the learning process, which can accumulate over time.

To address this, we propose two key modifications to ensure label consistency during augmentation. First, when replacing a masked feature, we restrict the replacement values only from samples of the same class. Although this setting was

considered in [76], we argue that preserving label consistency is even more critical in our iterative training framework. Second, after the initial training phase, we leverage the current model to predict the label of the augmented sample. The augmented sample is included in the dataset only if the model's prediction matches the original label to improve label consistency. However, it is important to note that feature consistency is not explicitly enforced during augmentation. This means the generated features may not always represent valid samples, potentially leading to inconsistencies.

### C. Confidence Calibration with Mixup

One of the key challenges with self-training is the issue of confirmation bias (also referred to as self-poisoning), where the model reinforces its incorrect predictions [11]. This occurs when the model makes an incorrect prediction, uses that prediction to update itself, and consequently becomes more confident in that incorrect prediction in future iterations.

Pseudo-labeling methods rely on thresholding to filter out low-confidence predictions and retain only high-confidence predictions. This process assumes that higher-confidence predictions are more likely to be correct. However, this assumption may not always hold due to issues with confidence calibration—the model's predicted confidence may not always reflect the true likelihood of the prediction being correct. By improving the model's confidence calibration, we can ensure that the pseudo-labels we incorporate into training are more accurate, thereby reducing the impact of self-poisoning.

To address this, we use *mixup* [81], a regularization technique to improve model generalization and calibration. Mixup generates new training samples by interpolating both features and labels between two randomly selected samples. Specifically, given two samples, the features and labels are mixed based on a mixing coefficient $\lambda$ drawn from a Beta distribution parameterized by $\alpha$. This results in fractional labels when mixing samples from different classes. As shown in [70], label interpolation is crucial for improving calibration. The target labels become fractional instead of one-hot, which changes the loss function for neural networks to the following:

$$\mathcal{L}_{\text{mixup}} = \lambda \cdot \mathcal{L}(y_i, \hat{y}_i) + (1 - \lambda) \cdot \mathcal{L}(y_j, \hat{y}_i)$$

where $\mathcal{L}$ is the binary cross-entropy loss, $y_i$ and $y_j$ are the true labels, and $\hat{y}_i$ is the model's prediction.

However, this approach is not applicable to models like Random Forests, which do not support fractional targets for classification tasks. To handle this, when mixing two samples with different labels, we assign the label corresponding to the sample with the higher interpolation coefficient. We provide evidence of how mixup improves confidence calibration in Section VII-D.

### V. DATASETS

We evaluate `ADAPT` across five publicly available malware detection datasets: two Android malware datasets, two Windows PE malware datasets, and one PDF malware dataset. Although the two Windows datasets share the same feature

TABLE I: Summary statistics of the datasets

| Dataset | Split | Duration | Benign Apps | Malicious Apps |
|---------|-------|----------|-------------|----------------|
| Drebin | Train | 2019-01 to 2019-12 | 21449 | 2187 |
|  | Validation | 2020-01 to 2020-06 | 10549 | 777 |
|  | Test | 2020-07 to 2021-12 | 20106 | 822 |
| APIGraph | Train | 2012-01 to 2012-12 | 16194 | 1447 |
|  | Validation | 2013-01 to 2013-06 | 13769 | 1446 |
|  | Test | 2013-07 to 2018-12 | 175048 | 13895 |
| BODMAS | Train | 2019-10 to 2019-12 | 13763 | 11082 |
|  | Validation | 2020-01 to 2020-03 | 13206 | 13769 |
|  | Test | 2020-04 to 2020-09 | 32608 | 27701 |
| EMBER | Train | 2018-01 | 29423 | 32040 |
|  | Validation | – | – | – |
|  | Test | 2018-02 to 2018-12 | 320577 | 358864 |
| PDF | Train | 10/13 to 10/27 | 25233 | 1208 |
|  | Validation | 10/27 to 09/10 | 70041 | 1893 |
|  | Test | 09/10 to 10/22 | 135751 | 20224 |

representation, they cover different time periods. The Android datasets differ both in feature sets and collection timelines. Table I summarizes key statistics for each dataset.

### A. Android Datasets

We use two Android malware detection datasets introduced by Chen et al. [24] for continuous learning using static features. The first dataset utilizes the Drebin feature set [12], while the second employs the APIGraph feature set [83].

The Drebin and APIGraph datasets represent Android applications using binary feature vectors that indicate the presence or absence of specific attributes. Drebin includes 16,978 dimensional features spanning eight categories: permissions, intents, API usage, and network addresses. In contrast, API-Graph reduces the feature space to 1,159 dimensions by clustering semantically similar APIs. Drebin samples were collected between 2019 and 2021, while APIGraph spans from 2012 to 2018. Both datasets were curated to mitigate experimental bias in malware detection [55], though they originally contained duplicates, which can affect result consistency [2], [84]. We, therefore, use deduplicated versions of both datasets in our experiments.

For both datasets, we use the first year's data as the training set, the next six months for validation, and the remaining months for testing. For pseudo-labeling, the training data is used to train the initial model, which is then periodically updated with data from subsequent months in the validation and test sets. The validation set is also used to jointly tune the model's hyperparameters and those specific to the pseudo-labeling algorithm.

### B. Windows Datasets

We use two Windows PE malware datasets: EMBER [7] and BODMAS [79]. Both datasets share identical feature sets but span different periods. Each file's features, extracted using the LIEF project [69], consist of byte sequences, imported functions, and header information compiled into a 2,381-dimensional feature vector.

Since both datasets have the same feature representation, we perform hyperparameter tuning only on the BODMAS dataset and use those tuned hyperparameters to evaluate the EMBER dataset. For the BODMAS dataset, we use the first three months of data for training, the next three months for validation, and the remaining six months for testing. For the EMBER dataset, we treat the first month's data as the labeled training set and use the remaining months for testing.

### C. PDF Dataset

We use the PDF malware dataset from [66] in our study. The dataset contains features representing benign and malicious PDFs over ten weeks, from 09/13/2012 to 10/22/2012. Although the number of features varies between weeks, we restrict our analysis to the features consistently present across all weeks to ensure compatibility with the model architectures used in our experiments. This results in 950-dimensional feature vectors. Compared to the Android and Windows datasets, this dataset exhibits less concept drift due to the shorter period.

For training, we randomly sample 10% of the data from the first two weeks. The following two weeks are used for validation, and the test set consists of data from the final six weeks.

## VI. EXPERIMENTAL SETTINGS

### A. Baseline Machine Learning Models

Our proposed algorithm is compatible with most machine learning models, requiring only access to model probabilities for pseudo-labeling, with data augmentation applied in a model-agnostic manner. We evaluate its performance using three widely adopted baseline models: Random Forest (RF) [21], XGBoost [23], and a Multi-Layer Perceptron (MLP) [60], chosen for their relevance in malware classification tasks [4], [7], [29], [55], [79], [83]. RF is computationally efficient, interpretable, and performs well on smaller datasets. XGBoost excels on tabular data, typical in malware detection, due to its ability to capture complex feature interactions [2], [36]. MLPs are well-suited for large-scale semi-supervised settings, supporting batch training and dynamic data augmentation, and exhibit robustness to noisy labels [76], making them ideal for pseudo-labeling methods [65].

### B. Self-Training Baselines

We compare our method against four self-training baselines: Adaptive Random Forest (ARF) [35], DroidEvolver++ (DE++) [41], Insomnia [9], and MORSE [76]. ARF and DE++ are online learning methods that address concept drift without requiring access to the original training dataset. In contrast, Insomnia and MORSE follow a semi-supervised paradigm, leveraging the original training data and pseudo-labeled samples for model updates. While our approach (ADAPT) also assumes access to the original training data, we additionally explore scenarios where this access is restricted. Further implementation details for all baseline methods are provided in Appendix A, and the restricted setting without original training data is discussed in Appendix D.

TABLE II: Comparison of methods across all datasets. Results are the mean of five runs (five different random seeds). Rows labeled **Oracle** represent an ideal scenario where retraining uses ground-truth labels; these are provided for reference only and are not considered in best result highlighting. For each dataset, <u>underlined</u> F1-scores indicate the best baseline (excluding proposed and Oracle methods), and **<u>bold+underline</u>** values indicate the best overall result (including proposed methods).

| Method | Model | Drebin | | | APIGraph | | | BODMAS | | | EMBER | | | PDF Malware | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | F1 | FPR | FNR | F1 | FPR | FNR | F1 | FPR | FNR | F1 | FPR | FNR | F1 | FPR | FNR |
| Offline ML | RF | 33.2 | 0.14 | 78.7 | 46.8 | 0.34 | 66.9 | 96.3 | 0.21 | 6.69 | 86.3 | 3.47 | 21.7 | <u>98.0</u> | 0.07 | 3.31 |
| | XGBoost | 44.9 | 0.66 | 65.6 | 66.5 | 0.97 | 42.7 | <u>99.2</u> | 0.30 | 1.23 | <u>92.3</u> | 3.58 | 11.4 | 97.9 | 0.14 | 3.17 |
| | MLP | 40.9 | 0.57 | 69.8 | 57.0 | 1.67 | 49.9 | 98.2 | 0.64 | 2.89 | 88.1 | 5.66 | 17.1 | 97.4 | 0.11 | 4.26 |
| Oracle | RF | 53.4 | 0.12 | 61.3 | 78.9 | 0.24 | 31.8 | 96.6 | 0.19 | 6.22 | 87.1 | 3.38 | 20.4 | 98.8 | 0.08 | 1.85 |
| | XGBoost | 68.3 | 0.52 | 41.2 | 85.5 | 0.67 | 17.6 | 99.6 | 0.21 | 0.61 | 93.1 | 3.46 | 10.1 | 99.2 | 0.09 | 1.04 |
| | MLP | 74.6 | 0.53 | 33.1 | 89.0 | 0.45 | 14.4 | 98.8 | 0.50 | 1.77 | 90.0 | 5.16 | 14.2 | 97.6 | 0.09 | 3.87 |
| Self-Training | ARF | 34.0 | 9.88 | 45.8 | 61.1 | 5.28 | 25.0 | 97.5 | 2.13 | 2.28 | 69.7 | 91.8 | 1.78 | 29.7 | 0.01 | 80.0 |
| | DE++ | 34.7 | 7.84 | 47.1 | 42.3 | 0.33 | 65.6 | 90.8 | 15.1 | 1.83 | 76.9 | 24.9 | 22.8 | 60.8 | 19.3 | 0.28 |
| | Insomnia | <u>49.8</u> | 0.70 | 60.6 | 68.1 | 2.05 | 31.8 | 97.2 | 0.56 | 4.81 | 87.1 | 6.06 | 18.5 | 88.2 | 0.09 | 15.4 |
| | MORSE | 41.2 | 4.50 | 54.7 | <u>75.3</u> | 1.82 | 22.2 | 98.8 | 0.48 | 1.73 | 89.8 | 6.34 | 13.7 | 97.4 | 0.12 | 4.16 |
| Proposed | RF+ADAPT | 49.8 | 0.21 | 64.3 | 67.2 | 0.34 | 46.0 | 96.1 | 0.22 | 7.11 | 86.2 | 3.56 | 21.7 | **<u>99.1</u>** | 0.09 | 1.26 |
| | XGBoost+ADAPT | **<u>57.1</u>** | 0.24 | 56.6 | 73.2 | 2.26 | 22.2 | **<u>99.6</u>** | 0.24 | 0.61 | **<u>93.5</u>** | 4.45 | 8.70 | 98.6 | 0.10 | 2.03 |
| | MLP+ADAPT | 50.5 | 0.77 | 59.1 | **<u>76.8</u>** | 0.81 | 29.4 | 98.6 | 0.41 | 2.26 | 89.1 | 5.01 | 16.0 | 97.2 | 0.09 | 4.73 |

### C. Hyperparameter Tuning

As shown in previous work [2], [24], continuous learning settings require dedicated hyperparameter tuning. In our experiments, we jointly tune the hyperparameters of both the baseline models and our proposed approach using the validation set. ADAPT introduces five key hyperparameters: the benign base threshold ($\tau_b$), the malware base threshold ($\tau_m$), the adaptive threshold weight ($\lambda$), the masking probability for augmentation ($p_a$), and the mixup parameter ($\alpha$). Each baseline self-training method also has its hyperparameters; for example, ARF uses a threshold ($\tau$), DE++ tunes the app buffer size, malware-to-benign ratio, and model aging threshold, Insomnia adjusts the pseudo-labeling fraction, and MORSE tunes the threshold along with weak and strong augmentation probabilities. A complete list of hyperparameters and their ranges is provided in Appendix G. To ensure fair comparison despite varying hyperparameter counts, we follow [2] and use a fixed budget of 200 random hyperparameter searches [16] for each method on each dataset.

### D. Evaluation

We use the same performance metrics as in [24]: the average F1-score, False Positive Rate (FPR), and False Negative Rate (FNR) across all months. Experiments are conducted using the optimal hyperparameters, with five random seeds on the test set. We report the mean and standard deviation of the performance metrics across these runs.

## VII. Results and Analysis

Table II presents the results of various methods across five datasets. XGBoost+ADAPT achieves the highest overall F1-score on the Drebin, BODMAS, and EMBER datasets, while MLP+ADAPT attains the best performance on the APIGraph dataset. Random Forest+ADAPT yields the top result on the PDF malware dataset. Across all benchmarks, ADAPT consistently enhances the performance of baseline machine learning

models, with particularly notable gains on the Android malware dataset, where concept drift causes greater performance degradation.

We also report results for an **Oracle** setting, where ground-truth labels are available for retraining, but augmentation and mixup components used in ADAPT are omitted. The results indicate that, for the Drebin and APIGraph datasets, the F1-score of ADAPT lags behind the Oracle model by more than 10%. However, this gap is less pronounced on the other three datasets, and on the EMBER dataset, XGBoost + ADAPT even surpasses the Oracle, likely due to the benefits introduced by data augmentation and mixup.

*Drebin Dataset*: Drebin is the most challenging dataset, as evidenced by the consistently low F1-scores of the offline machine learning models. Among the four self-training baselines, ARF, DE++, and MORSE do not yield notable improvements over the offline models in terms of F1-score. In contrast, Insomnia outperforms the offline MLP and XGBoost models, highlighting the effectiveness of our Insomnia adaptation for malware detection under concept drift. ADAPT achieves substantial improvements across all baseline models, with absolute F1-score gains of 16.6%, 12.2%, and 9.6% for RF, XGBoost, and MLP, respectively. This improvement is primarily attributed to a significant reduction in false negatives. Notably, our method does not increase the false positive rate; in fact, it reduces the false positive rate by 63.6% compared to the baseline XGBoost model.

Figure 6 illustrates the temporal dynamics of F1, FPR, and FNR over the test months for four models on the Drebin dataset. Notably, the XGBoost+ADAPT method consistently outperforms the other models in F1-score across most test months and is particularly effective at mitigating the sharp performance drop experienced by the baseline XGBoost model in certain months, such as month 9 (2021-03).

*APIGraph Dataset*: The baseline machine learning models achieve higher performance on the APIGraph dataset
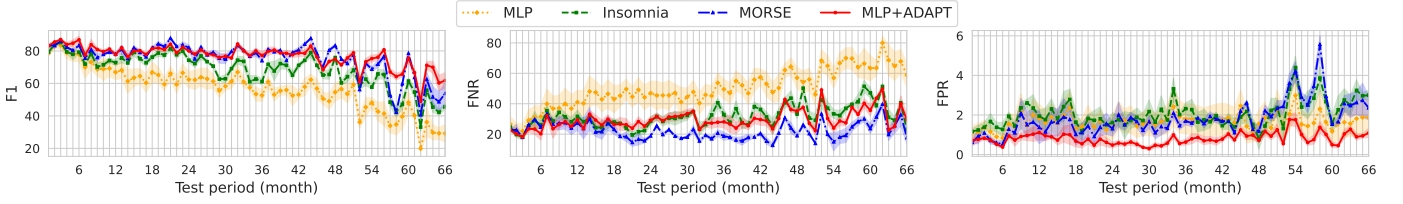
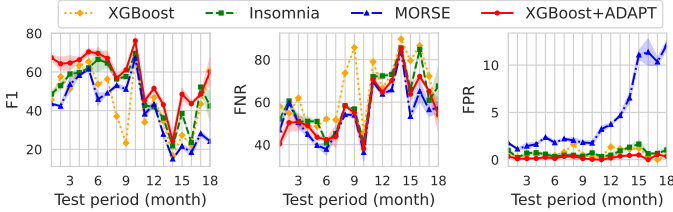Fig. 5: F1-score over test months on APIGraph dataset.



Fig. 6: Performance on the Drebin dataset: F1 score (left), FNR (middle), and FPR (right) over test months.
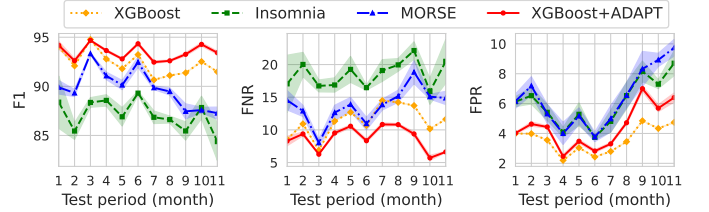


Fig. 8: Performance on the EMBER dataset: F1 score (left), FNR (middle), and FPR (right) over test months.
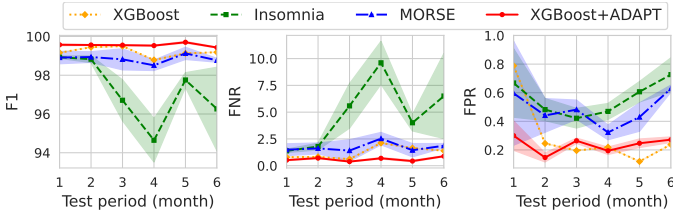


Fig. 7: Performance on the BODMAS dataset: F1 score (left), FNR (middle), and FPR (right) over test months.
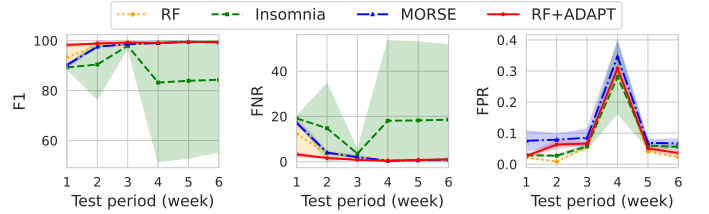


Fig. 9: Performance on the PDF dataset: F1 score (left), FNR (middle), and FPR (right) over test months.

compared to Drebin, as its features are specifically designed to be more resilient to concept drift [83]. Among the self-training baselines, all except DE++ show improvements over the offline models. Our proposed ADAPT method further improves upon the baselines, particularly for RF and MLP, both of which exhibit an approximately 20% increase in F1-score. While MORSE achieves the lowest False Negative Rate, the MLP+ADAPT method attains a lower False Positive Rate along with the best overall F1-score.

Figure 5 illustrates the F1-score over the test months for four models on the APIGraph dataset. The performance of all three adaptation models remains relatively stable during the first two years of test data, after which Insomnia experiences a sharp decline. ADAPT and MORSE perform similarly throughout most of the test period, except towards the end where more significant drift occurs.

*BODMAS Dataset:* The baseline models achieve notably high F1-scores on the BODMAS dataset, indicating that concept drift has a less pronounced effect in this setting. None of the self-training baseline methods surpasses the performance of the baseline XGBoost model, which already attains an F1-score of 99.2%. Our proposed ADAPT method further increases this score by 0.4%, while also reducing both the False Positive Rate (FPR) and False Negative Rate (FNR). This result underscores the value of our model-agnostic ap-

proach, as tree-based models such as XGBoost may inherently outperform neural networks on certain datasets [36]. A similar improvement is observed for MLP, although the performance of RF decreases. Figure 7 presents the monthly results for different models, highlighting that XGBoost+ADAPT consistently maintains high performance throughout the test period.

*EMBER Dataset:* Performance on the EMBER dataset is slightly lower, mainly due to the use of only one month of training data and the absence of separate hyperparameter tuning, in contrast to the three months of training used for BODMAS. On EMBER, the baseline XGBoost model outperforms all self-training methods, but our proposed ADAPT method still achieves a further 1.2% improvement in F1-score over the XGBoost baseline. Among the self-training baselines, DE++ and ARF perform significantly worse than the offline models, indicating greater susceptibility to self-poisoning and a high sensitivity to threshold settings in this dataset. Figure 8 shows the monthly performance of different models, highlighting that XGBoost+ADAPT consistently achieves the highest performance throughout the test period.

*PDF Dataset:* Figure 9 shows the weekly performance of different models on the PDF malware dataset. Random Forest outperforms both XGBoost and MLP, suggesting that more complex models may be susceptible to overfitting in this setting. Nonetheless, our proposed method, ADAPT, improves
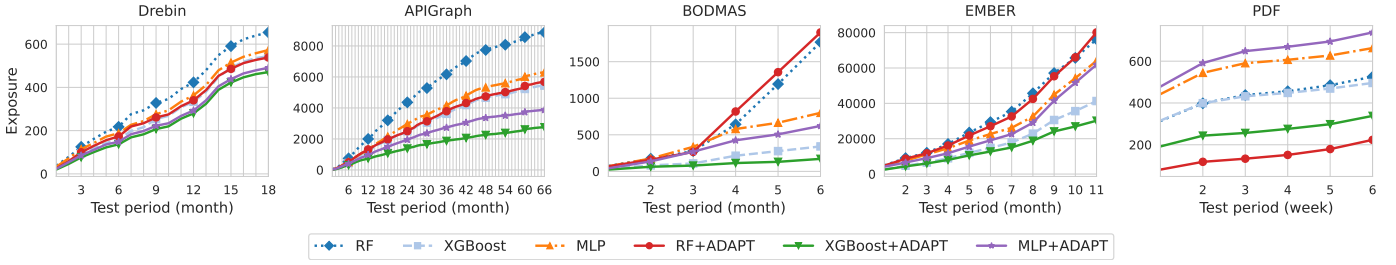
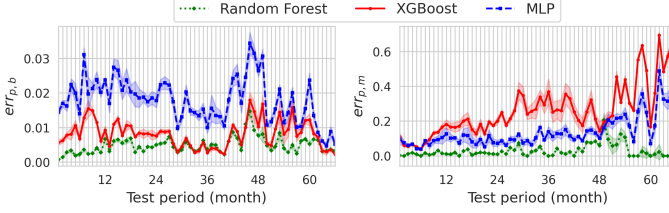Fig. 10: Exposure over test months on different datasets.



Fig. 11: Pseudo-labeling errors during the test months on the APIGraph dataset. (Left): Fraction of samples incorrectly pseudo-labeled as benign. (Right): Fraction of samples incorrectly pseudo-labeled as malware.

the F1-score by 1.1%, primarily through a reduction in the False Negative Rate (FNR).

### A. Exposure Analysis

To better capture the real-world risk associated with missed detections, we evaluate models using the Absolute Exposure (AE) metric introduced in [19]. AE quantifies the cumulative number of false negatives over time, reflecting the duration and extent to which users remain exposed to undetected threats. Unlike traditional accuracy metrics, AE directly measures how long malicious samples evade detection, providing a more practical assessment of security risk in streaming or temporally evolving scenarios.

Figure 10 presents the Absolute Exposure values for baseline and ADAPT models across all datasets. Notably, XGBoost+ADAPT achieves the lowest AE in every dataset except for the PDF malware dataset, where Random Forest + ADAPT attains the best performance.

### B. Effectiveness of ADAPT Across Models

We designed ADAPT as a model-agnostic adaptation strategy and demonstrated its ability to improve performance across different models. Our goal was to show that ADAPT enhances baseline models (RF, XGBoost, MLP) in the presence of concept drift. While different models may be preferable in various deployment settings, XGBoost+ADAPT emerges as a strong default choice, achieving the highest F1-score on three out of five datasets.

Although MLP+ADAPT achieves the highest performance over the full six-year test period on the APIGraph dataset, XGBoost+ADAPT (83.7% F1) outperforms both MLP+ADAPT (82.7% F1) and MORSE (80.5% F1) in the first year of test

data, a timeframe that more realistically represents practical deployment scenarios for unlabeled adaptation. This highlights that while ADAPT is model-agnostic, its strong synergy with XGBoost enhances its versatility as a broadly applicable solution. The statistical analysis of ADAPT, provided in Section VII-E further substantiates this finding.

### C. Error Analysis

Table III summarizes the changes in False Positive Rate (FPR) and False Negative Rate (FNR) between baseline models and those trained with ADAPT. Across datasets, F1-score improvements are mainly due to reduced FNR, except for Random Forest on BODMAS and MLP on PDF. This suggests that ADAPT enables the model to better learn from uncertain malware samples, improving its ability to correctly classify similar cases upon retraining. In most scenarios, FPR is maintained or reduced alongside FNR, though exceptions exist—e.g., XGBoost on APIGraph, which shows decreased FNR but increased FPR. Neural networks generally outperform other models in reducing FPR. We further examine this in the context of pseudo-labeling errors on APIGraph.

Figure 11 shows the pseudo-labeling errors introduced during the test months for the three ADAPT-trained models on the APIGraph dataset. Here, $err_{p,b}$ denotes the fraction of malware samples incorrectly pseudo-labeled as benign, and $err_{p,m}$ denotes benign samples mislabeled as malware. Across all models, $err_{p,b}$ is significantly lower than $err_{p,m}$, indicating that benign pseudo-labels are generally more accurate. This can be attributed to the higher thresholds for selecting benign samples: the optimal values of $\tau_b$ for Random Forest, XGBoost, and MLP are 0.96, 0.97, and 0.98, respectively. These high thresholds ensure that only samples with high benign confidence are selected, which helps prevent FNR degradation. This strategy is essential, as evasive malware may exhibit features similar to benign samples. Lowering the threshold could result in mislabeling such malware as benign, negatively impacting FNR.

Conversely, $err_{p,m}$ is notably higher due to the lower threshold $\tau_m$ used for pseudo-labeling malware samples—0.67, 0.70, and 0.63 for Random Forest, XGBoost, and MLP, respectively. These lower thresholds increase the likelihood of mislabeling benign samples as malware, a problem that intensifies over time, especially beyond four years. This explains the rise in FPR for XGBoost. Despite also introducing many noisy pseudo-labels, the MLP model reduces

TABLE III: Difference between False Positive Rate (FPR) and False Negative Rate (FNR) across different models and datasets. Positive values indicate a reduction in error, while negative values indicate an increase in error using the `ADAPT`.

| Model | Metric | Drebin | APIGraph | BODMAS | EMBER | PDF |
|---|---|---|---|---|---|---|
| RF | $\Delta$FPR | -0.07 | 0.00 | -0.01 | -0.09 | -0.02 |
| | $\Delta$FNR | 14.4 | 20.9 | -0.42 | 0.00 | 2.05 |
| XGBoost | $\Delta$FPR | 0.42 | -1.30 | 0.06 | -0.87 | 0.04 |
| | $\Delta$FNR | 9.00 | 20.5 | 0.62 | 2.70 | 1.14 |
| MLP | $\Delta$FPR | -0.20 | 0.86 | 0.23 | 0.65 | 0.02 |
| | $\Delta$FNR | 10.7 | 20.5 | 0.63 | 1.10 | -0.47 |

TABLE IV: Performance metrics (F1, FPR, FNR) for RF, XGBoost, and MLP with various components removed on the first six months of the APIGraph test set.

| Component Removed | RF | | | XGBoost | | | MLP | | |
|---|---|---|---|---|---|---|---|---|---|
| | F1 | FPR | FNR | F1 | FPR | FNR | F1 | FPR | FNR |
| None | 83.6 | 0.25 | 26.6 | 86.7 | 1.06 | 16.0 | 85.2 | 0.66 | 21.4 |
| Adaptive Th. | 80.8 | 0.17 | 31.2 | 86.4 | 1.00 | 16.8 | 81.2 | 0.50 | 25.1 |
| Augmentation | 82.3 | 0.21 | 28.8 | 85.6 | 1.11 | 17.5 | 84.8 | 0.75 | 21.4 |
| Mixup | 83.5 | 0.25 | 26.7 | 86.5 | 0.97 | 16.9 | 81.5 | 0.54 | 27.8 |

FPR compared to its baseline, likely due to the inherent robustness of neural networks to label noise in semi-supervised learning with data augmentation [76]. From a threat modeling perspective, using a lower threshold for malware selection is reasonable—benign apps lack the motivation to evade detection, so lower-confidence malware predictions still provide valuable drifted samples for retraining. While this increases the risk of labeling errors, it is essential for maintaining adaptability against evolving threats.

### D. Ablation and Mechanism Analysis

We conduct an ablation study by removing different components of `ADAPT` on the APIGraph dataset, as shown in Table IV for the first six months of test data. In this experiment, we replace the class-specific adaptive thresholding scheme with a fixed threshold of $\tau_b$. We observe that each component contributes to the overall improvement in F1-score and FNR, albeit varying degrees across different models. While replacing adaptive thresholding reduces the FPR marginally due to the higher threshold, as discussed in Section VII-C, it leads to a significant increase in the FNR, as drifted malware samples are not included in the retraining phase.

**Data Augmentation & Consistency.** Our feature-space augmentation does not guarantee valid samples, as modeling complex feature dependencies is intractable. To assess the impact, we analyze whether augmentation disrupts feature relationships by computing the mean absolute correlation in the original data and the Spearman correlation between original and augmented correlation matrices. As shown in Table V, low mean correlations indicate weak dependencies, and high Spearman values confirm that feature relationships are largely preserved, suggesting that the overall statistical structure remains stable.

TABLE V: Feature Correlations Before and After Augmentation

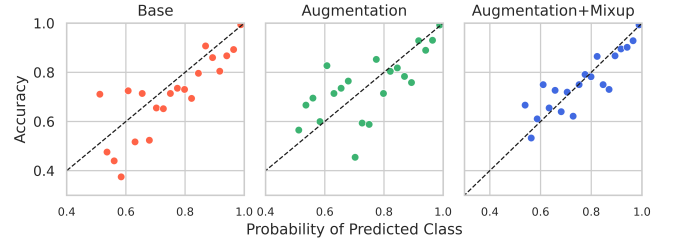| Dataset | Mean $|\rho(X)|$ | Spearman ($\rho(X)$, $\rho(X_{\text{aug}})$) |
|---|---|---|
| APIGraph | 0.0238 | 0.9247 |
| Drebin | 0.0269 | 0.9911 |
| Bodmas | 0.0276 | 0.9491 |
| PDF | 0.0857 | 0.9975 |



Fig. 12: Calibration plots for the XGBoost model on the APIGraph dataset. The model trained with Mixup shows better calibration, as the points are lying closer to the $y = x$ line.

**Effect of Mixup.** Figure 12 shows the effect of Mixup on calibration for XGBoost on APIGraph. The scatter plots compare accuracy versus confidence for three models: base (no augmentation or Mixup), with augmentation only, and with both augmentation and Mixup. The base model is notably overconfident, while data augmentation improves calibration, and combining it with Mixup provides the closest alignment between confidence and accuracy, especially for high-confidence predictions.

### E. Statistical Significance

Table VI presents the Wilcoxon signed-rank test p-values for comparing different machine learning models (Random Forest, MLP, and XGBoost) across multiple malware datasets. The Wilcoxon test [75] is a non-parametric statistical test used to assess whether two paired samples come from the same distribution. In this case, it evaluates whether the adapted model's F1-scores significantly differ from the baseline in the test months. Cells are shaded in gray if the p-value is less than 0.05, indicating statistically significant differences. Notably, we observe consistently low p-values for the XGBoost models across all datasets, suggesting a significant effect of adaptation on their performance.

TABLE VI: Wilcoxon $p$-values for different models across datasets. Cells with $p < 0.05$ are shaded.

| Dataset | Random Forest | MLP | XGBoost |
|---|---|---|---|
| Drebin | $7.63 \times 10^{-6}$ | $5.34 \times 10^{-5}$ | $1.53 \times 10^{-5}$ |
| APIGraph | $1.64 \times 10^{-12}$ | $1.64 \times 10^{-12}$ | $1.23 \times 10^{-10}$ |
| BODMAS | 0.6875 | 0.03125 | 0.03125 |
| EMBER | 0.7646 | 0.03223 | 0.00195 |
| PDF | 0.4375 | 0.15625 | 0.03125 |

TABLE VII: Performance on Drebin and APIGraph with active learning (AL) with 50 monthly annotation budget.

| Dataset | Method | F1 | FPR | FNR |
|---|---|---|---|---|
| Drebin | XGBoost | 62.5±1.95 | 0.57±0.04 | 47.7±1.85 |
|  | XGBoost + AL | 81.3±0.00 | 0.14±0.00 | 29.0±0.00 |
|  | HCC | 68.5±3.38 | 0.57±0.05 | 40.3±4.37 |
|  | XGBoost + `ADAPT` + AL | **81.7±0.28** | 0.16±0.01 | 27.8±0.44 |
| APIGraph | XGBoost | 79.2±0.43 | 0.98±0.06 | 24.2±1.20 |
|  | XGBoost + AL | 88.7±0.00 | 0.36±0.00 | 15.8±0.00 |
|  | HCC | 86.7±0.15 | 0.56±0.03 | 17.2±0.36 |
|  | XGBoost + `ADAPT` + AL | **89.4±0.17** | 0.79±0.02 | 9.55±0.18 |

## VIII. COMPATIBILITY WITH ACTIVE LEARNING

While `ADAPT` demonstrates significant improvements and can delay the onset of concept drift, at some point, labeled data may become necessary to handle drift effectively [41]. Active learning has proven to be one of the most effective approaches for addressing this issue in malware detection [24], [40]. Although the combination of semi-supervised and active learning has been explored in low-data scenarios [32], [71], its application to malware detection under concept drift remains largely unexplored. We explore this combination for the two Android malware detection datasets.

We follow the approach from [71], where the samples with the highest uncertainty (i.e., the lowest predicted probability) are selected for human annotation using active learning, while high-quality pseudo-labels are selected using `ADAPT` at every retraining step. We incorporate active learning into `ADAPT` with a slight modification in Algorithm 1. We adopt the same active learning setup as in [2], [24]. Specifically, we first sample the $k$ most uncertain samples for human annotation for each test month, where $k$ represents the monthly annotation budget. These samples are then added to the labeled training dataset $D_l$, which increases in size by $k$ each month. Following this, we perform Step 1 in Algorithm 1 to select pseudo-labeled samples using adaptive thresholding. However, we exclude the samples that have already been annotated through active learning, as we already have ground truth labels for them. We then continue with the rest of the algorithm and retrain the model. We perform a random hyperparameter search with 100 iterations, similar to [2].

Table VII presents the results with a monthly annotation budget of 50 on the Android datasets. For the XGBoost baseline model, we randomly select the same number of annotated samples per month to ensure consistency with the active learning methods following [10]. We report active learning results based on uncertainty sampling, using XGBoost and the Hierarchical Contrastive Learning Classifier (HCC) introduced in [24]. HCC, the state-of-the-art method for active learning with neural networks, leverages contrastive learning with a hierarchical loss function to distinguish between samples from different malware families in the embedding space.

From Table VII, we observe that active learning improves performance over random selection for the XGBoost model. Furthermore, `ADAPT` enhances this performance by leveraging the remaining unlabeled data. Specifically, we observe a 0.4%

TABLE VIII: Performance on multiclass classification (11-class) on BODMAS and EMBER datasets.

| Dataset | Method | F1 | Precision | Recall |
|---|---|---|---|---|
| BODMAS | XGBoost | 71.0±0.00 | 75.7±0.00 | 72.5±0.00 |
|  | ARF | 57.4±0.00 | 59.1±0.00 | 64.1±0.00 |
|  | MORSE | 68.1±2.21 | 67.8±4.68 | 73.1±0.78 |
|  | Insomnia | 70.1±1.21 | 71.6±0.99 | 73.6±1.70 |
|  | XGBoost+`ADAPT` | **71.4±0.65** | 75.6±1.27 | 72.6±0.61 |
| EMBER | XGBoost | 62.4±0.00 | 70.6±0.00 | 68.7±0.00 |
|  | ARF | 11.0±0.00 | 11.8±0.00 | 17.7±0.00 |
|  | MORSE | 58.4±1.49 | 68.4±1.11 | 64.5±1.02 |
|  | Insomnia | 52.7±0.78 | 66.4±0.28 | 57.7±0.67 |
|  | XGBoost+`ADAPT` | **65.2±0.65** | 70.3±0.94 | 71.4±0.59 |

TABLE IX: Performance on multiclass classification on BODMAS and EMBER datasets with active learning.

| Dataset | Method | F1 | Precision | Recall |
|---|---|---|---|---|
| BODMAS | XGBoost | 73.1±0.37 | 77.6±0.29 | 74.7±0.61 |
|  | XGBoost + AL | 73.9±0.00 | 79.7±0.00 | 75.4±0.00 |
|  | XGBoost+`ADAPT` + AL | **76.1±0.68** | 79.4±1.10 | 77.8±0.59 |
| EMBER | XGBoost | 78.0±0.51 | 80.6±0.35 | 82.4±0.54 |
|  | XGBoost + AL | 78.2±0.00 | 80.5±0.00 | 82.6±0.00 |
|  | XGBoost+`ADAPT` + AL | **80.6±1.00** | 81.8±1.10 | 84.5±0.68 |

increase in F1-score for the Drebin dataset and a 0.7% improvement for APIGraph. While these gains are relatively modest, they indicate that `ADAPT` can effectively benefit from active learning. Although our current approach involves a straightforward integration with active learning, exploring more sophisticated applications remains an avenue for future work.

## IX. EXTENSION TO MULTICLASS CLASSIFICATION

While we designed `ADAPT` primarily for binary malware detection, this section explores its applicability to a multiclass classification task. Specifically, we consider Windows malware family classification and construct two datasets from the BODMAS and EMBER datasets.

**Datasets.** We follow the same temporal split as in the binary classification task for training, validation, and test splits in both datasets. We select the top 10 most frequent malware families for each dataset in the training split and include the benign class, resulting in an 11-class classification task. We perform a hyperparameter search on the BODMAS dataset using 100 iterations per model and apply the selected hyperparameters directly to the EMBER dataset without further tuning. Appendix B provides additional details on these datasets.

**Algorithm Adjustments.** In the binary detection setting, we use separate threshold parameters, $\tau_b$ for benign samples and $\tau_m$ for malware samples, updated dynamically based on model predictions on unlabeled data using the adaptation parameter $\lambda$. While $\lambda$ is shared across all classes, introducing separate thresholds for each malware family in the multiclass setting would significantly increase the number of parameters, complicating hyperparameter tuning. To mitigate this, we share $\tau_m$ across all malware classes, keeping the number of hyperparameters the same as in the binary setting (five), regardless

of the number of malware families. Appendix B explores a scenario without a benign class, where only malware families are classified. Apart from this threshold adjustment, the rest of the algorithm, including augmentation and mixup, remains unchanged.

**Results.** We report the macro F1-score, precision, and recall in Table VIII for the XGBoost baseline, other self-training methods, and XGBoost+`ADAPT` on the two datasets. While `ADAPT` improves over the baseline on both datasets, other self-training methods perform worse than the XGBoost baseline. This effect is particularly pronounced in the EMBER dataset, which is more challenging and where we do not perform additional hyperparameter tuning. Notably, `ADAPT` improves the average F1-score over XGBoost by 2.8%, whereas other self-training methods perform significantly worse than the baseline, suggesting that `ADAPT` is more robust to hyperparameter values compared to other methods on EMBER feature.

**Unknown Family Detection.** To evaluate the open-set performance of the multi-class malware classifier, we conducted a monthly analysis with varying novelty in the unknown samples. For each month, we selected the top-$k$ most frequent unseen malware families ($k \in \{5, 10, 20\}$) and computed two key metrics: the *evasion success rate*—the fraction of unknown samples misclassified as benign—and the *AUC*, reflecting the model's ability to separate known from unknown samples based on confidence scores, in line with out-of-distribution (OOD) detection [78]. On the EMBER dataset, the XGBoost+`ADAPT` model achieved evasion success rates of 9.8%, 7.1%, and 6.6% for 5, 10, and 20 new families, respectively, with corresponding AUC values of 0.706, 0.713, and 0.713. These results highlight a tendency to misclassify unknown malware as known malicious classes, though the classifier's AUC remains stable as novelty increases.

**Active Learning.** We conduct additional experiments using an annotation budget of 50 samples per month for active learning with the XGBoost model. Table IX compares results for random sampling versus selecting the 50 most uncertain samples each month. Notably, integrating `ADAPT` with active learning improves performance over the random baseline by 3% on the BODMAS dataset and 2.6% on the EMBER dataset.

## X. LIMITATIONS & FUTURE WORK

**Evasion and Poisoning Attacks.** The effectiveness of our approach may be compromised by adversarial samples designed to evade malware detection classifiers [37], [73]. While we do not explicitly address adversarial examples, integrating adversarial robustness techniques [13] could enhance resilience. Additionally, poisoning attacks, where an adversary injects carefully crafted samples into the unlabeled data pool, can impact continual learning updates under concept drift [42], [67]. Such poisoning can cause the pseudo-labeling process to reinforce incorrect predictions, increasing the risk of persistent false negatives or false positives. Although semi-supervised learning methods may offer some resistance to random label noise [76], defending against targeted adversarial poisoning remains an open challenge [34]. We leave the exploration

of adversarial robustness and poisoning-aware defenses in pseudo-labeling pipelines as important directions for future work.

**Computational Efficiency.** `ADAPT` introduces no additional latency during deployment, as its inference time remains identical to that of the baseline model. However, the training time increases linearly due to data augmentation and mixup strategies. We provide a detailed analysis of the computational overhead in Appendix F.

**Limitations in Features.** Our experiments rely on static analysis features, which may not be sufficient for effectively handling concept drift [29]. Incorporating dynamic features, or a combination of static and dynamic features, could potentially improve the model's ability to adapt to concept drift. Its effectiveness with dynamic or hybrid features can be explored in future studies.

**Catastrophic forgetting.** Catastrophic forgetting refers to the tendency of a learned model to lose previously acquired knowledge when trained on new data—a challenge particularly relevant in malware detection [58]. We analyze its impact on the Android malware detection datasets in Appendix E. While our approach does not directly address catastrophic forgetting, future work could jointly integrate mechanisms to handle concept drift and mitigate forgetting. One promising direction is using replay-based methods from continual learning [59]. In our context, this could involve maintaining a replay buffer of class-balanced, confidently pseudo-labeled samples from intermediate time steps, which can be reused during subsequent model fine-tuning.

**Applicability to Other Security Domains.** Future work could explore extending `ADAPT` to other security domains, such as intrusion and phishing detection, where concept drift occurs frequently. While effective data augmentation may require domain-specific customization, other key components of `ADAPT`, such as adaptive thresholding and mixup regularization, remain domain-independent.

## XI. CONCLUSION

We present `ADAPT`, a pseudo-labeling-based approach to mitigate concept drift in malware detection. Our model-agnostic framework is validated across five diverse datasets exhibiting different types and degrees of distributional shift. Through comprehensive experiments, we show that `ADAPT` overcomes limitations of prior pseudo-labeling strategies and consistently improves detection performance. Our method also integrates well with active learning, further enhancing adaptability in practice. Overall, `ADAPT` provides a practical, label-efficient, and robust solution for maintaining detection accuracy under evolving threats, and its principles are readily extensible to other security and drift-sensitive tasks.

REFERENCES

[1] Hojjat Aghakhani, Fabio Gritti, Francesco Mecca, Martina Lindorfer, Stefano Ortolani, Davide Balzarotti, Giovanni Vigna, and Christopher Kruegel. When malware is packin'heat; limits of machine learning classifiers based on static analysis features. In *NDSS*, 2020.

[2] Md Tanvirul Alam, Dipkamal Bhusal, and Nidhi Rastogi. Revisiting static feature-based android malware detection. *arXiv preprint arXiv:2409.07397*, 2024.

[3] Md Tanvirul Alam, Romy Fieblinger, Ashim Mahara, and Nidhi Rastogi. Morph: Towards automated concept drift adaptation for malware detection. *arXiv preprint arXiv:2401.12790*, 2024.

[4] Mohammed S Alam and Son T Vuong. Random forest classification for detecting android malware. In *Proc. IEEE Int. Conf. Green Computing and Communications*, pages 663–669, 2013.

[5] David Alvarez-Melis and Nicolo Fusi. Geometric dataset distances via optimal transport. *Advances in Neural Information Processing Systems*, 33:21428–21439, 2020.

[6] Massih-Reza Amini, Vasilii Feofanov, Loic Pauletto, Lies Hadjadj, Emilie Devijver, and Yury Maximov. Self-training: A survey. *arXiv preprint arXiv:2202.12040*, 2022.

[7] H. S. Anderson and P. Roth. EMBER: An Open Dataset for Training Static PE Malware Machine Learning Models. *ArXiv e-prints*, April 2018.

[8] Hyrum S Anderson, Anant Kharkar, Bobby Filar, David Evans, and Phil Roth. Learning to evade static pe machine learning malware models via reinforcement learning. *arXiv preprint arXiv:1801.08917*, 2018.

[9] Giuseppina Andresini, Feargus Pendlebury, Fabio Pierazzi, Corrado Loglisci, Annalisa Appice, and Lorenzo Cavallaro. Insomnia: towards concept-drift robustness in network intrusion detection. In *Proceedings of the 14th ACM workshop on artificial intelligence and security*, 2021.

[10] Giovanni Apruzzese, Pavel Laskov, and Aliya Tastemirova. Sok: The impact of unlabelled data in cyberthreat detection. In *2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2022.

[11] Eric Arazo, Diego Ortego, Paul Albert, Noel E O'Connor, and Kevin McGuinness. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In *2020 International joint conference on neural networks (IJCNN)*, pages 1–8. IEEE, 2020.

[12] Daniel Arp, Michael Spreitzenbarth, Malte Hubner, Hugo Gascon, Konrad Rieck, and CERT Siemens. Drebin: Effective and explainable detection of android malware in your pocket. In *NDSS*, volume 14, pages 23–26, 2014.

[13] Tao Bai, Jinqi Luo, Jun Zhao, Bihan Wen, and Qian Wang. Recent advances in adversarial training for adversarial robustness. *arXiv preprint arXiv:2102.01356*, 2021.

[14] Federico Barbero, Feargus Pendlebury, Fabio Pierazzi, and Lorenzo Cavallaro. Transcending transcend: Revisiting malware classification in the presence of concept drift. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 805–823. IEEE, 2022.

[15] Anaël Beaugnon, Pierre Chifflier, and Francis R Bach. End-to-end active learning for computer security experts. In *AAAI Workshops*, pages 217–224, 2018.

[16] James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal of machine learning research*, 13(2), 2012.

[17] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in neural information processing systems*, 32, 2019.

[18] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.

[19] Marcus Botacin and Heitor Gomes. Towards more realistic evaluations: The impact of label delays in malware detection pipelines. *Computers & Security*, 148:104122, 2025.

[20] Tobias Braun, Irdin Pekaric, and Giovanni Apruzzese. Understanding the process of data labeling in cybersecurity. In *Proceedings of the 39th ACM/SIGAPP Symposium on Applied Computing*, pages 1596–1605, 2024.

[21] Leo Breiman. Random forests. *Machine learning*, 45:5–32, 2001.

[22] J Paul Brooks. Support vector machines with the ramp loss and the hard margin loss. *Operations research*, 59(2):467–479, 2011.

[23] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.

[24] Yizheng Chen, Zhoujie Ding, and David Wagner. Continuous learning for android malware detection. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 1127–1144, 2023.

[25] Yizheng Chen, Shiqi Wang, Dongdong She, and Suman Jana. On training robust pdf malware classifiers. In *Proceedings of the 29th USENIX Conference on Security Symposium*, pages 2343–2360, 2020.

[26] Zhi Chen, Zhenning Zhang, Zeliang Kan, Limin Yang, Jacopo Cortellazzi, Feargus Pendlebury, Fabio Pierazzi, Lorenzo Cavallaro, and Gang Wang. Is it overkill? analyzing feature-space concept drift in malware detectors. In *IEEE Security and Privacy Workshops (SPW)*. IEEE, 2023.

[27] Theo Chow, Zeliang Kan, Lorenz Linhardt, Lorenzo Cavallaro, Daniel Arp, and Fabio Pierazzi. Drift forensics of malware classifiers. In *Proceedings of the 16th ACM Workshop on Artificial Intelligence and Security*, pages 197–207, 2023.

[28] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 702–703, 2020.

[29] Savino Dambra, Yufei Han, Simone Aonzo, Platon Kotzias, Antonino Vitale, Juan Caballero, Davide Balzarotti, and Leyla Bilge. Decoding the secrets of machine learning in malware classification: A deep dive into datasets, feature extraction, and model performance. In *Proceedings of the 2023 ACM CCS*, 2023.

[30] Ido Finder, Eitam Sheetrit, and Nir Nissim. A time-interval-based active learning framework for enhanced pe malware acquisition and detection. *Computers & Security*, 121:102838, 2022.

[31] Cuiying Gao, Gaozhun Huang, Heng Li, Bang Wu, Yueming Wu, and Wei Yuan. A comprehensive study of learning-based android malware detectors under challenging environments. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, 2024.

[32] Mingfei Gao, Zizhao Zhang, Guo Yu, Sercan Ö Arık, Larry S. Davis, and Tomas Pfister. Consistency-based semi-supervised active learning: Towards minimizing labeling cost. In *Proc. ECCV 2020*, 2020.

[33] Josh Gardner, Zoran Popovic, and Ludwig Schmidt. Benchmarking distribution shift in tabular data with tableshift. *Advances in Neural Information Processing Systems*, 36, 2024.

[34] Micah Goldblum, Dimitris Tsipras, Chulin Xie, Xinyun Chen, Avi Schwarzschild, Dawn Song, Aleksander Madry, Bo Li, and Tom Goldstein. Dataset security for machine learning: Data poisoning, backdoor attacks, and defenses. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(2):1563–1580, 2022.

[35] Heitor M Gomes, Albert Bifet, Jesse Read, Jean Paul Barddal, Fabrício Enembreck, Bernhard Pfharinger, Geoff Holmes, and Talel Abdessalem. Adaptive random forests for evolving data stream classification. *Machine Learning*, 106:1469–1495, 2017.

[36] Léo Grinsztajn, Edouard Oyallon, and Gaël Varoquaux. Why do tree-based models still outperform deep learning on typical tabular data? *Advances in neural information processing systems*, 35:507–520, 2022.

[37] Kathrin Grosse, Nicolas Papernot, Praveen Manoharan, Michael Backes, and Patrick McDaniel. Adversarial perturbations against deep neural networks for malware classification. *arXiv preprint arXiv:1606.04435*, 2016.

[38] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.

[39] Maryam Issakhani, Princy Victor, Ali Tekeoglu, and Arash Habibi Lashkari. Pdf malware detection based on stacking learning. In *ICISSP*, pages 562–570, 2022.

[40] Roberto Jordaney, Kumar Sharad, Santanu K Dash, Zhi Wang, Davide Papini, Ilia Nouretdinov, and Lorenzo Cavallaro. Transcend: Detecting concept drift in malware classification models. In *26th USENIX security symposium (USENIX security 17)*, pages 625–642, 2017.

[41] Zeliang Kan, Feargus Pendlebury, Fabio Pierazzi, and Lorenzo Cavallaro. Investigating labelless drift adaptation for malware detection. In *Proceedings of the 14th ACM Workshop on Artificial Intelligence and Security*, pages 123–134, 2021.

[42] Łukasz Korycki and Bartosz Krawczyk. Adversarial concept drift detection under poisoning attacks for robust data stream mining. *Machine Learning*, 112(10):4013–4048, 2023.

[43] Ananya Kumar, Tengyu Ma, and Percy Liang. Understanding self-training for gradual domain adaptation. In *International Conference on Machine Learning*, pages 5468–5479. PMLR, 2020.

[44] Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3, 2013.

[45] Chi-Heng Lin, Chiraag Kaushik, Eva L Dyer, and Vidya Muthukumar. The good, the bad and the ugly sides of data augmentation: An implicit spectral regularization perspective. *Journal of Machine Learning Research*, 25(91):1–85, 2024.

[46] Martina Lindorfer, Matthias Neugschwandtner, and Christian Platzer. Marvin: Efficient and comprehensive mobile app classification through static and dynamic analysis. In *2015 IEEE 39th annual computer software and applications conference*. IEEE, 2015.

[47] Jie Lu, Anjin Liu, Fan Dong, Feng Gu, Joao Gama, and Guangquan Zhang. Learning under concept drift: A review. *IEEE transactions on knowledge and data engineering*, 31(12):2346–2363, 2018.

[48] Samaneh Mahdavifar, Andi Fitriah Abdul Kadir, Rasool Fatemi, Dima Alhadidi, and Ali A. Ghorbani. Dynamic android malware category classification using semi-supervised deep learning. In *Proc. IEEE DASC/PiCom/CBDCom/CyberSciTech*, pages 515–522, 2020.

[49] Brad Miller, Alex Kantchelian, Michael Carl Tschantz, Sadia Afroz, Rekha Bachwani, Riyaz Faizullabhoy, Ling Huang, Vaishaal Shankar, Tony Wu, George Yiu, et al. Reviewer integration and performance measurement for malware detection. In *DIMVA 2016, San Sebastián, Spain, July 7-8, 2016, Proceedings 13*. Springer, 2016.

[50] Annamalai Narayanan, Mahinthan Chandramohan, Lihui Chen, and Yang Liu. Context-aware, adaptive, and scalable android malware detection through online learning. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 1(3):157–175, 2017.

[51] Nir Nissim, Robert Moskovitch, Lior Rokach, and Yuval Elovici. Novel active learning methods for enhanced pc malware detection in windows os. *Expert Systems with Applications*, 41(13):5843–5857, 2014.

[52] Fakhroddin Noorbehbahani and Mohammad Saberi. Ransomware detection with semi-supervised learning. In *2020 10th International Conference on Computer and Knowledge Engineering (ICCKE)*, pages 024–029. IEEE, 2020.

[53] Rochester Institute of Technology. Research computing services, 2025. Accessed: 2025-07-31.

[54] Lucky Onwuzurike, Enrico Mariconti, Panagiotis Andriotis, Emiliano De Cristofaro, Gordon J. Ross, and Gianluca Stringhini. Mamadroid: Detecting android malware by building markov chains of behavioral models (extended version). *ACM Trans. Priv. Secur.*, 22(2), 2019.

[55] Feargus Pendlebury, Fabio Pierazzi, Roberto Jordaney, Johannes Kinder, Lorenzo Cavallaro, et al. Tesseract: Eliminating experimental bias in malware classification across space and time. In *Proceedings of the 28th USENIX Security Symposium*. USENIX Association, 2019.

[56] Joaquin Quiñonero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset shift in machine learning*. Mit Press, 2022.

[57] Dima Rabadi and Sin G Teo. Advanced windows methods on malware detection and classification. In *Proceedings of the 36th Annual Computer Security Applications Conference*, pages 54–68, 2020.

[58] Mohammad Saidur Rahman, Scott Coull, and Matthew Wright. On the limitations of continual learning for malware classification. In *Conference on Lifelong Learning Agents*, pages 564–582. PMLR, 2022.

[59] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. *Advances in neural information processing systems*, 32, 2019.

[60] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088), 1986.

[61] Igor Santos, Javier Nieves, and Pablo G Bringas. Semi-supervised learning for unknown malware detection. In *International Symposium on Distributed Computing and Artificial Intelligence*, pages 415–422. Springer, 2011.

[62] Weiwei Shi, Yihong Gong, Chris Ding, Zhiheng MaXiaoyu Tao, and Nanning Zheng. Transductive semi-supervised deep learning using minmax features. In *ECCV*, pages 299–315, 2018.

[63] Methaq A Shyaa, Noor Farizah Ibrahim, Zurinahni Zainol, Rosni Abdullah, Mohammed Anbar, and Laith Alzubaidi. Evolving cybersecurity frontiers: A comprehensive survey on concept drift and feature dynamics aware machine and deep learning in intrusion detection systems. *Engineering Applications of Artificial Intelligence*, 137:109143, 2024.

[64] Rami Sihwail, Khairuddin Omar, and KA Zainol Ariffin. A survey on malware analysis techniques: Static, dynamic, hybrid and memory analysis. *Int. J. Adv. Sci. Eng. Inf. Technol*, 8(4-2):1662–1671, 2018.

[65] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33:596–608, 2020.

[66] Nedim Šrndić and Pavel Laskov. Hidost: a static machine-learning-based detector of malicious files. *EURASIP Journal on Information Security*, 2016:1–20, 2016.

[67] Rahim Taheri, Reza Javidan, Mohammad Shojafar, Zahra Pooranian, Ali Miri, and Mauro Conti. On defending against label flipping attacks on malware detection systems. *Neural Computing and Applications*, 32:14781–14800, 2020.

[68] Kimberly Tam, Ali Feizollah, Nor Badrul Anuar, Rosli Salleh, and Lorenzo Cavallaro. The evolution of android malware and android analysis techniques. *ACM Computing Surveys (CSUR)*, 49(4), 2017.

[69] Romain Thomas. Lief - library to instrument executable formats. https://lief.quarkslab.com/, apr 2017.

[70] Sunil Thulasidasan, Gopinath Chennupati, Jeff A Bilmes, Tanmoy Bhattacharya, and Sarah Michalak. On mixup training: Improved calibration and predictive uncertainty for deep neural networks. *Advances in neural information processing systems*, 32, 2019.

[71] Katrin Tomanek and Udo Hahn. Semi-supervised active learning for sequence labeling. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 1039–1047, 2009.

[72] Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2008.

[73] 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*, 26(2):930–966, 2023.

[74] Yanshuo Wang, Jie Hong, Ali Cheraghian, Shafin Rahman, David Ahmedt-Aristizabal, Lars Petersson, and Mehrtash Harandi. Continual test-time domain adaptation via dynamic sample selection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1701–1710, 2024.

[75] Frank Wilcoxon. Individual comparisons by ranking methods. In *Breakthroughs in statistics: Methodology and distribution*, pages 196–202. Springer, 1992.

[76] Xian Wu, Wenbo Guo, Jia Yan, Baris Coskun, and Xinyu Xing. From grim reality to practical solution: Malware classification in real-world noise. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 2602–2619. IEEE Computer Society, 2023.

[77] Ke Xu, Yingjiu Li, Robert Deng, Kai Chen, and Jiayun Xu. Droidevolver: Self-evolving android malware detection system. In *2019 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2019.

[78] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. *International Journal of Computer Vision*, 132(12):5635–5662, 2024.

[79] Limin Yang, Arridhana Ciptadi, Ihar Laziuk, Ali Ahmadzadeh, and Gang Wang. Bodmas: An open dataset for learning based temporal analysis of pe malware. In *2021 IEEE Security and Privacy Workshops (SPW)*, pages 78–84. IEEE, 2021.

[80] Limin Yang, Wenbo Guo, Qingying Hao, Arridhana Ciptadi, Ali Ahmadzadeh, Xinyu Xing, and Gang Wang. Cade: Detecting and explaining concept drift samples for security applications. In *USENIX security symposium*, pages 2327–2344, 2021.

[81] Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.

[82] Linjun Zhang, Zhun Deng, Kenji Kawaguchi, Amirata Ghorbani, and James Zou. How does mixup help with robustness and generalization? *arXiv preprint arXiv:2010.04819*, 2020.

[83] Xiaohan Zhang, Yuan Zhang, Ming Zhong, Daizong Ding, Yinzhi Cao, Yukun Zhang, Mi Zhang, and Min Yang. Enhancing state-of-the-art classifiers with api semantics to detect evolved android malware. In *Proceedings of the 2020 ACM SIGSAC conference on computer and communications security*, pages 757–770, 2020.

[84] Yanjie Zhao, Li Li, Haoyu Wang, Haipeng Cai, Tegawendé F Bissyandé, Jacques Klein, and John Grundy. On the impact of sample duplication in machine-learning-based android malware detection. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 30(3):1–38, 2021.

DroidEvolver [77] and its improved version, DroidEvolver++ (DE++) [41], are notable self-training approaches for Android malware detection, but both suffer from self-poisoning and limited evaluation across datasets. Beyond these, research on self-training for concept drift in malware detection remains scarce, motivating our inclusion of additional established self-training algorithms for comparison.

1) **Adaptive Random Forest (ARF) [35]**: ARF is an online ensemble method for concept drift adaptation. To enable self-training, we introduce a threshold parameter for pseudo-label selection, which is tuned in conjunction with other tree hyperparameters, allowing ARF to update using pseudo-labeled samples without requiring ground-truth labels.

2) **DroidEvolver++ (DE++) [41]**: DE++ improves upon DroidEvolver [77] for malware drift adaptation without requiring ongoing ground truth. It uses an ensemble of five linear models to generate pseudo-labels and tracks model aging with a buffer of recent samples. When prediction consistency in the buffer drops below a threshold, the model is updated using new pseudo-labels from the ensemble.

3) **Insomnia [9]**: Insomnia adapts to concept drift in network intrusion detection via co-training between a Nearest Centroid (NC) classifier and an MLP, with the NC providing pseudo-labels for uncertain MLP predictions. Both models are updated using pseudo-labeled data. For malware detection, we found NC to be ineffective, so we use XGBoost with the MLP, and select the most confident samples for updating to avoid early self-poisoning.

4) **MORSE [76]**: MORSE is a state-of-the-art semi-supervised method for malware family classification with noisy labels. It uses the FixMatch [65] framework, incorporating an augmentation strategy that enforces consistency between weakly and strongly augmented versions of each pseudo-labeled sample. We include MORSE as a baseline due to its strong performance in the presence of label noise.

The first two baselines, ARF and DE++, are online methods that update one sample at a time and do not require access to the full training dataset after initialization. In contrast, Insomnia and MORSE use batch training, updating models by combining original training data with pseudo-labeled samples. While our main experiments assume the training data remains available, we also consider restricted scenarios without access to the original training data (see Appendix D).

A. *Dataset Statistics*

Table X and Table XI present class-wise statistics for the BODMAS and EMBER datasets, respectively.

TABLE X: Statistics of the BODMAS multiclass classification dataset

| ID | Family | Train Samples | Validation Samples | Test Samples |
|---|---|---|---|---|
| 0 | Benign | 3000 | 3000 | 6000 |
| 1 | Wacatac | 2337 | 917 | 988 |
| 2 | Upatre | 930 | 1148 | 1369 |
| 3 | Mira | 586 | 512 | 791 |
| 4 | Small | 558 | 1026 | 1475 |
| 5 | Dinwod | 405 | 682 | 701 |
| 6 | Wabot | 368 | 748 | 2405 |
| 7 | Autorun | 363 | 148 | 60 |
| 8 | Musecador | 309 | 717 | 19 |
| 9 | Gepys | 298 | 382 | 255 |
| 10 | Berbew | 284 | 195 | 1195 |

TABLE XI: Statistics of the EMBER multiclass classification dataset

| ID | Family | Train Samples | Test Samples |
|---|---|---|---|
| 0 | Benign | 500 | 5500 |
| 1 | InstallMonster | 500 | 3217 |
| 2 | AdPoshel | 500 | 2647 |
| 3 | Zusy | 500 | 3783 |
| 4 | Fareit | 500 | 4662 |
| 5 | Emotet | 500 | 5386 |
| 6 | DealPly | 500 | 2227 |
| 7 | DotDo | 500 | 953 |
| 8 | StartSurf | 500 | 3819 |
| 9 | Mira | 425 | 304 |
| 10 | Tiggre | 405 | 1284 |

For BODMAS, we utilize all available samples per malware family in training and randomly sample 1,000 benign samples per month across all splits. Some families, like Gepys, are rare in the test set, while others, such as Berbew, are more frequent during testing than training. For EMBER, we cap each class, including Benign, at 500 samples in training and per month in testing to manage dataset size. This produces approximately balanced training data but class imbalance in the test set.

B. *Experiments Excluding Benign Class*

In practice, malware analysis typically proceeds in two steps: first, binary classification distinguishes between benign and malicious samples; second, if a sample is identified as malware, its family is determined [58]. We emulate this scenario for malware family classification by removing benign samples from the BODMAS (Table X) and EMBER (Table XI) datasets, resulting in a 10-class task. In this setting, ADAPT uses a single threshold parameter $\tau_m$ shared across all malware families; otherwise, the algorithm is unchanged from the binary case. This reduces the number of ADAPT hyperparameters from five to four.

Table XII presents the results of ADAPT and other self-training approaches across both datasets. ADAPT outperforms the baseline XGBoost methods on both datasets, while other self-training methods perform worse. The improvement is especially notable on the EMBER dataset, where ADAPT achieves a 6.5% increase in the F1-score. Table XIII presents

TABLE XII: Performance on multiclass classification (10-class) on BODMAS and EMBER datasets.

| Dataset | Method | F1 | Precision | Recall |
|---------|--------|-----|-----------|--------|
| BODMAS | XGBoost | 76.8±0.00 | 82.9±0.00 | 78.3±0.00 |
| | ARF | 63.5±0.00 | 70.6±0.00 | 64.0±0.00 |
| | MORSE | 75.3±0.86 | 80.1±1.18 | 76.4±1.09 |
| | Insomnia | 62.2±1.19 | 69.5±1.47 | 65.7±1.44 |
| | XGBoost+ADAPT | **77.0±1.07** | 82.2±1.26 | 77.7±0.99 |
| EMBER | XGBoost | 58.7±0.00 | 67.8±0.00 | 66.3±0.00 |
| | ARF | 10.9±0.00 | 9.66±0.00 | 21.1±0.00 |
| | MORSE | 53.6±2.23 | 62.8±1.68 | 63.1±0.89 |
| | Insomnia | 39.6±0.53 | 58.5±1.31 | 49.0±0.23 |
| | XGBoost+ADAPT | **65.2±1.43** | 75.1±1.94 | 70.4±1.35 |

TABLE XIII: Performance on multiclass classification (10-class) on BODMAS and EMBER datasets with active learning.

| Dataset | Method | F1 | Precision | Recall |
|---------|--------|-----|-----------|--------|
| BODMAS | XGBoost | 81.2±0.96 | 86.2±0.99 | 81.9±0.96 |
| | XGBoost + AL | 83.2±0.00 | 89.2±0.00 | 83.6±0.00 |
| | XGBoost+ADAPT + AL | **85.5±1.08** | 90.6±0.37 | 85.8±1.58 |
| EMBER | XGBoost | 76.4±1.22 | 79.2±1.06 | 80.6±1.21 |
| | XGBoost + AL | 77.1±0.00 | 79.6±0.00 | 81.8±0.00 |
| | XGBoost+ADAPT + AL | **83.8±0.65** | 85.0±0.55 | 87.2±0.40 |

the results for active learning with a monthly annotation budget of 50 samples. Active learning with uncertainty sampling outperforms the baseline XGBoost model trained with the same annotation budget using random sampling. Additionally, incorporating ADAPT further enhances performance on both datasets. Specifically, integrating ADAPT improves the F1-score by 2.3% on the BODMAS dataset and 6.1% on the EMBER dataset, demonstrating the effectiveness of leveraging unlabeled samples even in the presence of active learning.

## APPENDIX C
### SUPPORT VECTOR MACHINE(SVM) EXPERIMENTS

We conduct additional experiments using SVM baseline models on two Android datasets, APIGraph and Drebin, using the same hyperparameter selection strategy from Section VI-C. The results, including those with ADAPT, are summarized in Table XIV.

The baseline SVM model performs competitively with the baseline XGBoost model on the APIGraph dataset. However, its performance is significantly lower on the Drebin dataset than the XGBoost baseline.

Integrating ADAPT with the SVM model proves to be effective, yielding performance improvements over the baseline SVM across both datasets. Notably, ADAPT significantly reduces the false negative rate, improving the F1-score. This effect is particularly pronounced in the Drebin dataset, where we observe an average F1-score increase of 14.2%.

## APPENDIX D
### SOURCE-FREE ADAPTATION

In some scenarios, access to the original training data may be unavailable after initial model training due to privacy, proprietary, storage, or computational constraints [77]. In these

TABLE XIV: SVM performance on APIGraph and Drebin (average over five runs).

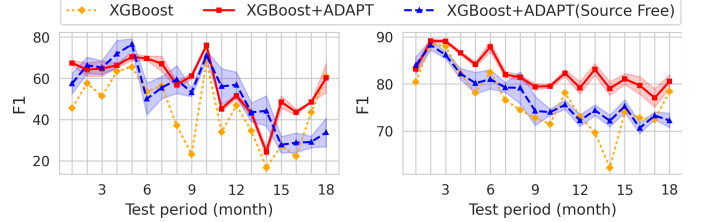| Method | APIGraph | | | Drebin | | |
|--------|----------|-----|-----|--------|-----|-----|
| | F1 | FPR | FNR | F1 | FPR | FNR |
| SVM | 66.3 | 1.41 | 39.5 | 33.8 | 0.70 | 76.7 |
| SVM+ADAPT | 73.6 | 1.47 | 29.3 | 48.0 | 0.84 | 61.5 |



Fig. 13: F1 score on Drebin (left) and APIGraph (right) datasets across the first 18 test months for XGBoost, XGBoost with ADAPT, and XGBoost with ADAPT in a source-free setting (i.e., without access to the training data during retraining).

cases, model updates must rely solely on pseudo-labeled data. We explore this setting for XGBoost on the Drebin and APIGraph datasets, modifying Algorithm 1 at line 7 so that $D_m = D_p$, meaning only pseudo-labeled data are available after the initial phase. Separate hyperparameter optimization is performed for this restricted scenario using the same search space.

Figure 13 compares the baseline XGBoost model, the ADAPT semi-supervised update, and source-free adaptation on these datasets. For APIGraph, results are shown for the first 18 test months, as performance declines over longer periods (average F1: 61.24%), though the source-free method offers noticeable gains in this window. On Drebin, source-free adaptation improves the average F1 by 11.3% over the baseline in the first year, lagging the semi-supervised variant by just 1.76%. On APIGraph, it outperforms the baseline by 0.82% but is 3.97% below the semi-supervised approach. These findings indicate that ADAPT can still mitigate concept drift when the original training data is unavailable, making it suitable for continual learning scenarios with storage limitations [58].

## APPENDIX E
### CATASTROPHIC FORGETTING

Catastrophic forgetting, where a model loses previously acquired knowledge when learning from new data, is a significant challenge in continual learning [58]. While our work does not directly address this issue, we assess its impact on the Android malware datasets using XGBoost and MLP. Specifically, we evaluate the final model trained on the last month ($M_T$) against the first year of test data, comparing its F1 score, FPR, and FNR to those of models trained during each monthly retraining phase.

The results are presented in Figure 14 for the Drebin dataset and Figure 15 for the APIGraph dataset. Both models exhibit some degree of catastrophic forgetting on the Drebin dataset,
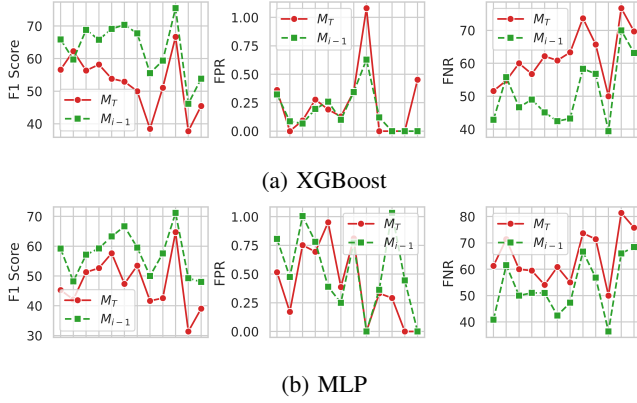
Fig. 14: Impact of catastrophic forgetting on Drebin: XGBoost and MLP performance over the first 12 test months. $M_T$ is the final-month model; $M_{i-1}$ is the previous month's model.
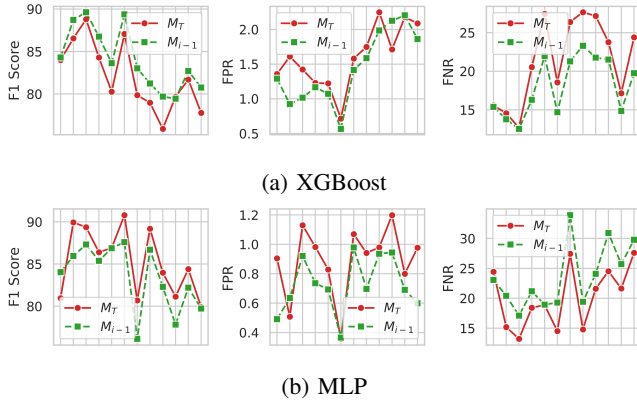


Fig. 15: Impact of catastrophic forgetting on the APIGraph dataset for XGBoost and MLP models.

as the final model's performance degrades compared to earlier models. In contrast, on the APIGraph dataset, the final MLP model outperforms its earlier versions, indicating no signs of catastrophic forgetting. For the XGBoost model, catastrophic forgetting is primarily driven by an increase in false negatives, whereas for the MLP model, the FNR increases on Drebin but decreases on APIGraph.

## APPENDIX F
### COMPUTATIONAL EFFICIENCY

While ADAPT does not increase inference time, it incurs additional training cost, which increases approximately linearly due to data augmentation and mixup. This overhead is primarily driven by two factors: (1) the time required to generate augmented samples and (2) the additional time needed to train on the expanded dataset, roughly three times larger than the original training set.

The extent of this increase depends on the dataset size, feature dimensionality, and model complexity. Table XV reports the training time overhead introduced by ADAPT across different models and datasets. The most significant increase is

TABLE XV: Training time (in seconds) of different baseline models with ADAPT

| Dataset | RF | RF + ADAPT | MLP | MLP + ADAPT | XGBoost | XGBoost + ADAPT |
|---------|------|-------|-------|--------|---------|----------|
| Drebin | 38.5 | 281.0 | 158.0 | 712.9 | 90.3 | 1378.6 |
| APIGraph | 19.5 | 105.0 | 32.0 | 37.0 | 6.2 | 17.2 |
| BODMAS | 45.9 | 188.3 | 242.9 | 681.6 | 93.5 | 301.6 |
| EMBER | 111.7 | 422.1 | 512.7 | 1618.5 | 237.4 | 668.1 |
| PDF | 12.5 | 57.5 | 38.1 | 70.4 | 6.1 | 66.0 |

observed for the XGBoost model on the Drebin dataset, where the high feature dimensionality (16,978) substantially impacts computation time. For other datasets, training time typically increases by three to five relative to the baseline model.

## APPENDIX G
### HYPERPARAMETER SEARCH

Table XVI summarizes the hyperparameter ranges used for each method. For self-training baselines, we adopt default parameters for the five online learning classifiers. Insomnia, MORSE, and our method ADAPT use different underlying models: Insomnia applies co-training with XGBoost and MLP, MORSE uses MLP, and ADAPT is evaluated with Random Forest, XGBoost, and MLP. All models share the same parameter search space as in the offline setting. We perform a random search with 200 trials across the joint parameter space to select the best hyperparameters for each dataset.

## APPENDIX H
### THEORETICAL ANALYSIS

Our self-training approach for malware detection under concept drift is theoretically supported by prior work on gradual self-training [43]. We outline key theoretical results explaining its effectiveness under distribution shifts.

### A. Setup and Assumptions

We consider a gradual domain shift scenario, characterized by a sequence of distributions $P_0, P_1, \ldots, P_T$, where $P_0$ is the source domain and $P_T$ the target. The shift is assumed gradual: for a small $\epsilon > 0$, the distance $\rho(P_t, P_{t+1}) < \epsilon$ for all $0 \leq t < T$, with $\rho$ denoting a distributional distance metric. We assume there exists a sequence of classifiers $\theta_t$ that can correctly classify most samples with a margin, and that the incremental shifts between domains permit effective iterative adaptation.

### B. Main Theoretical Results

The primary result from [43] states that under gradual domain shift, self-training can iteratively reduce error. Specifically:

*Theorem 1 ( [43], Theorem 3.2):*

Let $P$ and $Q$ be two distributions such that the Wasserstein distance [72] satisfies $\rho(P, Q) = \rho < 1/R$, where $R$ is the regularization parameter. Assume that both distributions have the same marginal distribution over labels, i.e., $P(Y) = Q(Y)$.

Suppose we have an initial model $\theta \in \Theta_R$ trained on $P$, and we perform self-training using $n$ unlabeled samples from

TABLE XVI: Hyperparameter search spaces for various methods. "U" indicates sampling from a uniform random distribution within the specified range. Hyperparameters marked with "*" are only used during the retraining phase. The Insomnia method also uses hyperparameters from MLP and XGBoost, while MORSE uses those from MLP. ADAPT utilizes the baseline model, which can be either Random Forest (RF), XGBoost, or MLP.

| Model | Hyperparameter | Candidate Values |
|---|---|---|
| Random Forest | n_estimators | $2^x$, where $x \in U[5, 10]$ |
| | max_depth | $2^y$, where $y \in U[5, 10]$ |
| | criterion | {gini, entropy, log_loss} |
| | class_weight | {None, "balanced"} |
| XGBoost | max_depth | $2^w$, where $w \in U[3, 7]$ |
| | alpha | $10^a$, where $a \in U[-8, 0]$ |
| | lambda | $10^b$, where $b \in U[-8, 0]$ |
| | eta | $3.0 \times 10^c$, where $c \in U[-2, -1]$ |
| | balance | {True, False} |
| | num_boost_round | {100, 150, 200, 300, 400} |
| MLP | mlp_layers | {[100, 100], [512, 256, 128], [512, 384, 256, 128], [512, 384, 256, 128, 64]} |
| | learning_rate | $10^d$, where $d \in U[-5, -3]$ |
| | dropout | $x$, where $x \in U[0.0, 0.5]$ |
| | batch_size | $2^e$, where $e \in \{5, 6, 7, 8, 9, 10\}$ |
| | epochs | {25, 30, 35, 40, 50, 60, 80, 100, 150} |
| | optimizer | {Adam} |
| | balance | {True, False} |
| | cont_learning_epochs* | {0.1, 0.2, 0.3, 0.4, 0.5} |
| SVM | C | $10^z$, where $z \in U[-4, 3]$ |
| | class_weight | {None, "balanced"} |
| ARF | n_models | $2^x$, where $x \in U[3, 5]$ |
| | max_features | {sqrt, log2, None} |
| | max_depth | $2^y$, where $y \in U[5, 10]$ |
| | lambda_value | 6 |
| | threshold | $x$, where $x \in U[0.6, 0.99]$ |
| DE++ | age_threshold_low | $x$, where $x \in U[0.0, 0.5]$ |
| | buffer_ratio | $x$, where $x \in U[0.1, 0.5]$ |
| | buffer_size | {1000, 2000, 3000} |
| Insomnia | fraction | $x$, where $x \in U[0.0, 0.8]$ |
| MORSE | threshold | $x$, where $x \in U[0.8, 0.99]$ |
| | weak_augment | $x$, where $x \in U[0.0, 0.1]$ |
| | strong_augment | $x$, where $x \in U[0.1, 0.2]$ |
| ADAPT | threshold_benign | $x$, where $x \in U[0.8, 0.99]$ |
| | threshold_malware | $x$, where $x \in U[0.6, 0.99]$ |
| | lambda | $x$, where $x \in U[0.0, 0.5]$ |
| | mask_ratio | $x$, where $x \in U[0.0, 0.2]$ |
| | mixup_alpha | $x$, where $x \in U[0.0, 0.2]$ |

$Q$ to obtain an updated model $\theta'$. Then, with probability at least $1 - \delta$ over the sampling of $n$ unlabeled examples from $Q$, the ramp loss [22] of the new model $\theta'$ on $Q$ satisfies:

$$L_r(\theta', Q) \leq \frac{2}{1 - \rho R} L_r(\theta, P) + \alpha^*$$
$$+ \frac{4BR + \sqrt{2 \log(2/\delta)}}{\sqrt{n}}, \quad (6)$$

Here, $L_r(\theta, P)$ represents the ramp loss of the initial model on the source distribution $P$, while $\alpha^* = \min_{\theta^* \in \Theta_R} L_r(\theta^*, Q)$ denotes the minimum achievable ramp loss on the target distribution $Q$ within the hypothesis space $\Theta_R$. The term $B$ is an upper bound on the norm of the input features, ensuring bounded feature magnitudes. The parameter $R$ controls the capacity of the hypothesis space by regularizing the complexity of the learned models. Finally, $\delta$ is the confidence parameter

that determines the probability with which the bound holds over the sampling of $n$ unlabeled examples from $Q$.

This theorem implies that if the initial model has low loss on the source domain and the distributional shift is gradual, self-training will yield a model with controlled error on the new distribution. Applying this argument iteratively across $T$ steps leads to an exponential improvement over direct adaptation:

*Corollary 1 ( [43], Corollary 3.3):* Under the assumptions of $\alpha^*$-separation, no label shift, gradual shift, and bounded data, suppose the initial model $\theta_0$ has a ramp loss of at most $\alpha_0 \geq \alpha^*$ on the source distribution $P_0$, i.e., $L_r(\theta_0, P_0) \leq \alpha_0$. If self-training is applied iteratively over $T$ steps to produce the final model $\theta_T$, denoted as $\theta_T = ST(\theta_0, (S_1, \ldots, S_T))$, then with high probability, the ramp loss on the final target distribution $P_T$ is bounded as:

$$L_r(\theta_T, P_T) \leq \beta^{T+1}\left(\alpha_0 + \frac{4BR + \sqrt{2\log(2T/\delta)}}{\sqrt{n}}\right)$$
(7)

where $\beta = \frac{2}{1-\rho R}$. This result demonstrates that gradual self-training effectively controls the error over multiple adaptation steps, preventing the potential failure cases that can arise from direct adaptation.

### C. Extension to ADAPT

In the malware classification setting, we extend the theoretical framework by considering different distribution shifts for the two classes: malware and benign. Let $\rho_m$ and $\rho_b$ denote the Wasserstein shifts for the malware and benign classes, respectively, where we assume $\rho_m > \rho_b$. This reflects the practical scenario where malware samples undergo more significant distributional changes compared to benign samples.

*1) Asymmetric Error:* Following Lemma A.2 from [43], we analyze class-specific error bounds under different shift magnitudes. The classification error on the target distribution $Q$ for malware and benign classes is bounded by:

$$\text{Err}^m(\theta, Q) \leq \frac{1}{1 - \rho_m R} L_r^m(\theta, P),$$
$$\text{Err}^b(\theta, Q) \leq \frac{1}{1 - \rho_b R} L_r^b(\theta, P).$$
(8)

Since $\rho_m > \rho_b$, the multiplicative factor for malware is larger, meaning that—even with similar initial losses—classification error for malware will be higher than for benign samples under shift.

Consequently, malware samples experience greater drift and are more likely to be misclassified as benign, increasing the false negative rate (FNR). By contrast, benign samples cross into the malware region less frequently, resulting in a lower false positive rate (FPR). This asymmetry explains why FNR typically exceeds FPR in malware classification under distributional shift.

*2) Impact of Class-Specific Thresholding:* To mitigate the imbalance between false negative and false positive rates, we introduce separate thresholds $\tau_m$ and $\tau_b$ for malware and benign classifications, respectively. Instead of classifying all samples, we only classify those for which the model's confidence exceeds the respective threshold. This effectively reduces the number of classified samples but also lowers the classification error.

Formally, let $Q_{\tau_m}^m$ and $Q_{\tau_b}^b$ denote the distributions of malware and benign samples that remain after thresholding. Since thresholding removes uncertain samples near the decision boundary, the ramp loss on the remaining classified samples decreases:

$$L_r^m(\theta, Q_{\tau_m}^m) \leq L_r^m(\theta, Q),$$
$$L_r^b(\theta, Q_{\tau_b}^b) \leq L_r^b(\theta, Q).$$
(9)

Applying the error bound from Equation (8) to the thresholded samples, we obtain:

$$\text{Err}^m(\theta, Q) \leq \frac{1}{1 - \rho_m R} L_r^m(\theta, P_{\tau_m}),$$
$$\text{Err}^b(\theta, Q) \leq \frac{1}{1 - \rho_b R} L_r^b(\theta, P_{\tau_b}).$$
(10)

Since thresholding removes samples that are close to the decision boundary—where errors are more likely—the filtered distributions $Q_{\tau_m}^m$ and $Q_{\tau_b}^b$ contain fewer misclassified samples, leading to lower classification errors.

Thresholding reduces the number of classified samples, which affects the generalization bound in Equation (6). Specifically, the third term in the bound, which depends on the sample size $n$, is modified as follows:

$$\frac{4BR + \sqrt{2\log(2/\delta)}}{\sqrt{n}} \quad \rightarrow \quad \frac{4BR + \sqrt{2\log(2/\delta)}}{\sqrt{n_{\tau_m,\tau_b}}},$$
(11)

where $n_{\tau_m,\tau_b}$ is the number of samples retained after thresholding at $\tau_m$ and $\tau_b$, and we have $n_{\tau_m,\tau_b} < n$. Since the denominator decreases, this term increases, leading to a looser bound.

*3) Modified Error Bound:* Following Lemma A.3 from [43], we account for the fact that the total error depends on the mixture of malware and benign samples in $Q$. Given that $q_m$ and $q_b$ represent the proportions of malware and benign samples in $Q$ (i.e., $q_m + q_b = 1$), we replace the expectation in the bound with an explicit weighted sum.

First, the error term in Equation (6) becomes:

$$\frac{2}{1 - \rho R} L_r(\theta, P) \quad \rightarrow \quad q_m \frac{2}{1 - \rho_m R} L_r^m(\theta, P_{\tau_m})$$
$$+ q_b \frac{2}{1 - \rho_b R} L_r^b(\theta, P_{\tau_b}).$$
(12)

Thus, the final modified bound becomes:

$$L_r(\theta', Q_{\tau_m,\tau_b}) \leq q_m \frac{2}{1 - \rho_m R} L_r^m(\theta, P_{\tau_m})$$
$$+ q_b \frac{2}{1 - \rho_b R} L_r^b(\theta, P_{\tau_b})$$
$$+ \alpha^*$$
$$+ \frac{4BR + \sqrt{2\log(2/\delta)}}{\sqrt{n_{\tau_m,\tau_b}}}.$$
(13)

*4) Impact of Augmentation & Mixup:* Data augmentation and mixup serve as effective regularizers during training [45], [82], helping to reduce overfitting and improve generalization. In our theoretical framework, this regularization implicitly reduces the norm of the learned weight vector $w$, leading to a smaller capacity parameter $R$ (i.e., $\|w\|_2 \leq R$). A lower $R$ constrains model complexity, making it less likely to fit noise or learn overly intricate decision boundaries. Consequently, stronger regularization tightens the bound in Equation 13 and enhances robustness to distributional shifts. A detailed theoretical analysis of these effects is left for future work.