# LAMDA: A LONGITUDINAL ANDROID MALWARE BENCHMARK FOR CONCEPT DRIFT ANALYSIS

**Anonymous authors** 

Paper under double-blind review

## **ABSTRACT**

Machine learning (ML)-based malware detection systems often fail to account for the dynamic nature of real-world training and test data distributions. In practice, these distributions evolve due to frequent changes in the Android ecosystem, adversarial development of new malware families, and the continuous emergence of both benign and malicious applications. Prior studies have shown that such concept drift—distributional shifts in benign and malicious samples, leads to significant degradation in detection performance over time. Despite the practical importance of this issue, existing datasets are often outdated and limited in temporal scope, diversity of malware families, and sample scale, making them insufficient for the systematic evaluation of concept drift in malware detection.

To address this gap, we present LAMDA, the largest and most temporally diverse Android malware benchmark to date, designed specifically for concept drift analysis. LAMDA spans 12 years (2013–2025, excluding 2015), includes over 1 million samples (approximately 37% labeled as malware), and covers 1,380 malware families and 150,000 singleton samples, reflecting the natural distribution and evolution of real-world Android applications. We empirically demonstrate LAMDA's utility by quantifying performance degradation of standard ML models over time and analyzing feature stability across years. As the most comprehensive Android malware dataset to date, LAMDA enables in-depth research into temporal drift, generalization, explainability, and evolving detection challenges. The dataset and code are available at: https://anonymous.4open.science/r/LAMDA-D763/.

# 1 Introduction

Android malware poses a growing threat to user privacy and security, with over 33 million attacks blocked in 2024 alone (Kaspersky, 2024; AV-TEST, 2025). Static feature based ML methods, which analyze features extracted from Android application packages (APKs), have emerged as a promising defense mechanism (Arp et al., 2014; Mariconti et al., 2017). However, these detectors often suffer performance degradation over time due to *concept drift* — gradual shifts in the feature distribution caused by the evolving nature of both malicious and benign software (Yang et al., 2021d).

Concept drift can result from several factors, including changes in developer practices, updates to Android APIs, and, most significantly, the evolving and adaptive strategies of malware authors (Greenberg, 2020). To evade detection, adversaries frequently obfuscate or modify their code by injecting alternative API calls, altering manifest components, or exploiting newly introduced services (News, 2024). For example, the Android trojan *SoumniBot* obfuscates its manifest file to evade analysis and detection (News, 2024). These tactics lead to observable shifts in static features over time, undermining the robustness of ML-based detection systems (Yang et al., 2021d). Prior studies have shown that malware families (i.e., clusters of samples exhibiting similar behavioral traits) play a central role in driving such drifts (Chow et al., 2023; Barbero et al., 2022b).

Although concept drift plays a central role in Android malware evolution, most existing datasets are not designed to support drift analysis. Datasets such as Drebin (Arp et al., 2014), TESSER-ACT (Pendlebury et al., 2019), and APIGraph (Zhang et al., 2020) are limited in temporal coverage, family diversity, or structural organization for studying drift. Similarly, Windows-based datasets like EMBER (Anderson & Roth, 2018b), SOREL-20M (Harang & Rudd, 2020), and BODMAS (Yang et al., 2021c) are constrained by short collection periods or target different ecosystems. While EM-

BERSim (Corlatescu et al., 2023), MalNet (Freitas et al., 2020), and AnoShift (Dragoi et al., 2022) offer task-specific contributions, they do not support longitudinal drift analysis in Android malware classification. To address these gaps, we introduce LAMDA, a novel Android malware benchmark dataset curated for temporal drift analysis with family evolution. LAMDA spans over 12 years (i.e., 2013–2025, excluding 2015 due to the unavailability of hashes in the AndroZoo repository (Allix et al., 2016)), covering 1,008,381 APK samples across 1,380 unique malware families and over 150,000 Singleton samples (i.e., samples without *av class* labels) from AndroZoo repository (Alecci et al., 2024). Each sample is labeled using VirusTotal's vt\_detection count (VirusTotal, 2025) reported in AndroZoo database (Alecci et al., 2024). The samples are decompiled to extract finegrained static features following the Drebin (Arp et al., 2014) feature definitions.

We validate LAMDA through a series of comprehensive evaluations, including longitudinal degradation analysis of the supervised binary classification under concept drift (AnoShift-style (Dragoi et al., 2022)), temporally disjoint training (testing), and family-wise feature stability assessments. LAMDA enables explanation-guided analysis of concept drift and combines long-term structural modeling with SHAP-based attributions (Lundberg & Lee, 2017), allowing researchers to trace how feature relevance shifts over time and better understand the underlying causes of model degradation.

In summary, the contributions of this paper are as follows:

- We present LAMDA, a large-scale Android malware benchmark comprising over 1 million APKs across 1,380 unique families spanning for 12 years (2013 to 2025, excluding 2015), built on static features based on Drebin (Arp et al., 2014) features.
- We perform a detail concept drift detection using structured temporal splits (Dragoi et al., 2022) to show that LAMDA exhibits pronounced distributional shift than prior benchmark.
- We conduct comprehensive drift analysis, including per-feature distribution shifts, feature stability analysis across malware families (Zhang et al., 2020), temporal label flipping analysis, and SHAP-based explanation (Lundberg & Lee, 2017) drift that reveal temporal changes in feature importance.
- We show that existing drift adaptation methods, though effective on prior benchmark (Zhang et al., 2020), fail to generalize on LAMDA due to its more realistic and pronounced concept drift.

## 2 Related Work

In this section, we discuss prior work and their limitations that motivate the creation of LAMDA.

Evolution of Malware Datasets and Benchmarks. Early malware datasets such as Drebin (Arp et al., 2014) (Android) and EMBER (Anderson & Roth, 2018b) (Windows) have played a pivotal role to study concept drift in malware analysis. More recent efforts—including SOREL-20M (Harang & Rudd, 2020) and BODMAS (Yang et al., 2021c) for Windows, and TESSERACT (Pendlebury et al., 2019), APIGraph (Zhang et al., 2020), and Chen et al. (2023) for Android attempt to address limitations in scale and recency. Nonetheless, these datasets suffer from one or more major limitations—they are often outdated, contain either relatively few malware samples or families, or lack long-term temporal coverage necessary for studying the evolution of malware. For example, Drebin spans only 2010–2012 with 5,560 samples from 179 families; TESSERACT covers 2014–2016 with 12,735 samples; API Graph spans 2012–2018 with 32,089 samples from 1,120 families; and Chen et al. (2023) includes 10,200 samples across 254 families from 2019–2021. Despite their temporal spread, these datasets are not explicitly structured to support longitudinal drift analysis or capture evolutionary patterns in malware behavior.

**Explainability and Semantic Features.** Explainability is critical for understanding how feature importance shifts under concept drift. While Drebin (Arp et al., 2014) and BODMAS (Yang et al., 2021b) introduced interpretable features and temporal structure, few studies have systematically used them to analyze drift. TRANSCENDENT (Barbero et al., 2022a) incorporates semantic reasoning for selective prediction, but longitudinal robustness of explanations remains underexplored due to limited dataset support. LAMDA fills this gap by providing a temporally structured benchmark with interpretable features and SHapley Additive exPlanations (SHAP)-based explanations (Lundberg & Lee, 2017), enabling fine-grained, longitudinal analysis of model behavior and drift.

# 3 LAMDA CREATION

 In this section, we describe the construction process of the LAMDA. We have downloaded APKs from AndrooZoo repository (Allix et al., 2016) and decompiled APKs to extract static Drebin (Arp et al., 2014) features and then transformed the features into binary vectors for downstream analysis.

**Label Assignment and Collection Strategy.** To construct a large-scale, temporally diverse dataset, we use metadata from AndroZoo (Allix et al., 2016), including APK hashes, VirusTotal (VT) results, and submission dates. For each year from 2013 to 2025 (excluding 2015, which lacks valid entries), we collect APKs and assign binary labels using the vt\_detection field. Following prior heuristics (Arp et al., 2014; Pendlebury et al., 2019), we define: (i) *Benign* for vt\_detection = 0, (ii) *Malware* for vt\_detection  $\geq$  4, and (iii) discard *Uncertain* samples with scores in 1, 3. The  $\geq$  4 threshold mitigates label noise by requiring stronger AV consensus (Chen et al., 2016).

To reduce sampling bias in learning systems, we collected 50,000 malware and 50,000 benign samples per year, while preserving month-wise temporal distributions across both categories. Although prior work such as TESSERACT (Pendlebury et al., 2019; Chen et al., 2023) adopts a 90:10 benign-to-malware ratio, we attempt to maintain a balanced 50:50 ratio (Anderson & Roth, 2018b). This choice is motivated by the need to mitigate the risk of skewed learned representations (such as overfitting (Shwartz-Ziv et al., 2023), disparity in learning (Zhou et al., 2023)) that can arise from class imbalance. A balanced dataset helps ensure that the model learns meaningful distinctions between classes, captures a wider range of malware families, and is exposed to a broader spectrum of behaviors and evasive techniques. Such diversity not only enables longitudinal generalization studies but also increases the difficulty of the detection task, particularly for learning systems that must contend with rare, novel, or semantically similar malware families (Anderson & Roth, 2018b). Nonetheless, due to limited availability of malware samples in certain years such as 2017, 2023, 2024, and 2025, LAMDA still exhibits class imbalance, with each of these years showing different imbalance ratios.

Another practical challenge during data collection involved download and decompilation failures, requiring us to over-fetch APKs to meet target counts. To mitigate this, we included a 20% overhead in the number of APK hashes per year. All APKs are retrieved via authenticated academic access to the AndroZoo repository <sup>1</sup> and stored in a consistent directory structure ([year]/malware/, [year]/benign/) to facilitate temporal slicing and cross-year analysis. Corrupted or undecompilable samples are excluded and logged for transparency. The final dataset comprises over one million APKs. A detailed year-wise breakdown is provided in Appendix A.

Family Label Acquisition. To enable finer-grained analysis of how malware behavior evolves over time, we assign family-level labels to all malware samples using AVClass2 (Sebastián et al., 2016), which standardizes noisy antivirus vendor labels into consistent malware family names. This is critical for developing detection systems that generalize to emerging threats. The labeling process involves retrieving VirusTotal (2025) reports, converting them to the required format, running AVClass2, and post-processing the output to retain SHA256 mappings. Figures 1(a) and 1(b) show the yearly distribution of recurring versus newly observed families and the count of singleton families—those that appear only once—respectively. Family labels enable research into more complex tasks such as multi-class classification and the study of temporal trends across malware families.

Decompilation and Static Feature Extraction. Each APK is statically decompiled using apktool (Brut, 2025), producing a disassembled smali representation and the original AndroidManifest.xml. We parse these artifacts to extract a diverse set of static features commonly used in Android malware detection (Arp et al., 2014). Specifically, the AndroidManifest.xml file is analyzed to obtain the list of requested permissions (e.g., ACCESS\_FINE\_LOCATION), declared activities and services, broadcast receivers, required hardware components, and intent filters (Arp et al., 2014). Meanwhile, the disassembled smali code is scanned to identify invocations of restricted APIs (e.g., NotificationManager.notify), suspicious API usages (e.g., getSystemService), and embedded hardcoded IPs/URLs (e.g., e.crashlytics.com). The extracted Drebin feature sets comprises several static categories derived from Android APKs (Arp et al., 2014). A detailed list of features is provided in Appendix B.

https://androzoo.uni.lu/access



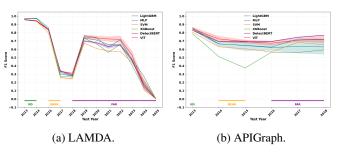


Figure 1: Temporal trends in malware family evolution.

Figure 2: F1-score over time across different temporal splits.

**Vectorization and Temporal Feature Alignment.** After decompiling each APK, we extract static features into a .data file. Each year's data is split into 80% training and 20% testing sets using stratified sampling to preserve class balance. From the training set, we construct a global vocabulary by taking the union of unique tokens across all samples, yielding 9,690,482 (≈ 9.69 million) raw features (Yang et al., 2021d). Each APK is then represented as a high-dimensional binary vector using a bag-of-tokens model, where each token corresponds to a binary feature indicating its presence or absence in the sample (Arp et al., 2014). To reduce dimensionality and ensure computational feasibility, we apply VarianceThreshold from scikit-learn to eliminate lowvariance features. For all experiments, we use the Baseline variant, which applies a threshold of 0.001 (Rahman et al., 2025), resulting in 4,561 final features. This compact and consistent representation supports a range of downstream tasks, including supervised learning, drift analysis, and continual learning. More details are in Appendix B. The dataset is initially created in a sparse matrix format, storing binary feature vectors and metadata as compressed .npz files to optimize for storage and computational efficiency. These . npz files are organized by year and stratified into training and test splits. A detailed breakdown of feature dimensions under varying VarianceThreshold settings is provided in Appendix B. For scalability of LAMDA, we have also published global features, variance threshold objects and selected features after applying VarianceThreshold.

## 4 CONCEPT DRIFT DETECTION

In this section, we examine performance degradation of supervised models across temporally distant splits (Dragoi et al., 2022) (Section 4.1) to detect concept drift, followed by distributional shifts using Jeffreys divergence and t-SNE visualizations (Sections 4.1 and 4.2).

# 4.1 CONCEPT DRIFT DETECTION WITH SUPERVISED LEARNING

**Experimental Setting.** To evaluate the robustness of malware detectors under temporal distribution shifts, we perform supervised learning experiments using four widely adapted detector models from the malware research — Linear SVM, LightGBM, MLP, XGBoost, detectBERT, and ViT (Arp et al., 2014; Anderson & Roth, 2018b). Detailed model configurations are provided in Appendix D.

Following the AnoShift benchmark (Dragoi et al., 2022), we divide LAMDA into three temporally separated regions: TRAIN (i.e., initial training set), IID, NEAR, and FAR. TRAIN+IID includes samples from 2013–2014, with the last month of each year held out for IID evaluation. Models are trained on all other months, and the held-out portion serves as an in-distribution test set to measure baseline performance on temporally adjacent, unseen data. To assess generalization under drift, we define two evaluation regions: NEAR (2016–2017) and FAR (2018–2025), allowing systematic analysis of performance degradation as the temporal gap from training increases. For comparison, we evaluate the same models on the APIGraph dataset (Zhang et al., 2020) using a similar split: training on 2012, IID on 2013, NEAR on 2014, and FAR on 2015–2018. All experiments are repeated five times with different random seeds. For each split, we report results in Table 1 as mean±std, averaged across all runs and years within each split. Figure 2 shows yearly results averaged over runs. This experiment uses the LAMDA-baseline with a VarianceThreshold of 0.001. To study the impact of feature space on drift, we also test thresholds of 0.01 and 0.0001, with results provided in Appendix D.

Table 1: Comparison of performances on LAMDA and API Graph across three temporal splits.

Split	Model			LAMDA					API Graph		
Spite	Model	F1	ROC-AUC	PR-AUC	FNR	FPR	F1	ROC-AUC	PR-AUC	FNR	FPR
***	LightGBM MLP	97.49 <sub>±0.17</sub> 97.21 <sub>±0.12</sub>	99.55 <sub>±0.03</sub> 99.48 <sub>±0.04</sub>	99.50 <sub>±0.11</sub> 99.38 <sub>±0.20</sub>	1.74 <sub>±0.34</sub> 2.50 <sub>±0.06</sub>	2.69 <sub>±0.48</sub> 2.58 <sub>±0.85</sub>	85.95 <sub>±0.00</sub> 85.79 <sub>±0.00</sub>	98.91 <sub>±0.00</sub> 96.37 <sub>±0.00</sub>	95.20 <sub>±0.00</sub> 88.49 <sub>±0.00</sub>	$22.39_{\pm 0.00}$ $20.31_{\pm 0.00}$	0.33 <sub>±0.00</sub> 0.67 <sub>±0.00</sub>
IID	SVM	$94.98_{\pm 1.07}$	$98.89_{\pm0.28}$	$98.75_{\pm 0.46}$	$4.82_{\pm 0.76}$	$4.09_{\pm 0.55}$	$82.00_{\pm0.00}$	$97.33_{\pm0.00}$	$90.94_{\pm0.00}$	$26.74_{\pm0.00}$	$0.60_{\pm 0.00}$
	XGBoost	$97.05_{\pm0.14}$	$99.15_{\pm0.16}$	$97.68_{\pm 1.16}$	$2.20_{\pm 0.43}$	$2.96_{\pm0.01}$	$80.33_{\pm0.00}$	$96.05_{\pm0.00}$	89.74+0.00	$28.35_{\pm0.00}$	$0.75_{\pm 0.00}$
	DetectBERT	$95.27_{\pm 0.86}$	$98.92_{\pm 0.27}$	$98.61_{\pm 0.79}$	$3.96_{\pm0.13}$	$4.48_{\pm 0.59}$	$83.05_{\pm0.00}$	$98.82_{\pm0.00}$	$94.19_{\pm 0.00}$	$25.60_{\pm0.00}$	$0.48_{\pm 0.00}$
	ViT	$94.97_{\pm 1.59}$	$98.91_{\pm 0.37}$	$98.61_{\pm 0.85}$	$4.63_{\pm0.81}$	$4.17_{\pm 0.73}$	86.64 <sub>±0.00</sub>	$98.65_{\pm0.00}$	$93.73_{\pm 0.00}$	$17.82_{\pm 0.00}$	$0.81_{\pm 0.00}$
	LightGBM	$59.48_{\pm 28.20}$	$74.05_{\pm 23.76}$	70.18 $_{\pm 27.10}$	$50.51_{\pm 30.82}$	$1.85_{\pm 0.95}$	66.77 <sub>±0.00</sub>	$95.94_{\pm0.00}$	$83.01_{\pm 0.00}$	$47.68_{\pm0.00}$	$0.48_{\pm 0.00}$
NEAR	MLP	$56.57_{\pm 28.41}$	$82.71_{\pm 11.19}$	$67.94_{\pm 24.59}$	$51.95_{\pm 30.42}$	$3.98_{\pm 1.19}$	$68.72_{\pm 0.00}$	$86.79_{\pm0.00}$	$68.31_{\pm 0.00}$	$38.80_{\pm 0.00}$	$1.87_{\pm 0.00}$
NLAK	SVM	$52.91_{\pm 28.40}$	$75.18_{\pm 17.98}$	$62.53_{\pm 29.82}$	$55.62_{\pm 28.88}$	$4.71_{\pm 0.97}$	$63.06_{\pm0.00}$	$89.25_{\pm0.00}$	$70.15_{\pm0.00}$	$45.48_{\pm0.00}$	$2.03_{\pm 0.00}$
	XGBoost	$55.84_{\pm 29.73}$	$77.75_{\pm 16.85}$	$68.14_{\pm 26.55}$	$53.84_{\pm 30.94}$	$2.14_{\pm0.86}$	$51.47_{\pm 0.00}$	$86.90_{\pm0.00}$	$65.69_{\pm0.00}$	$60.41_{\pm 0.00}$	$1.57_{\pm 0.00}$
	DetectBERT	$59.11_{\pm 34.70}$	$73.86_{\pm 29.73}$	$66.95_{\pm 37.60}$	$49.92_{\pm 37.55}$	$4.05_{\pm0.84}$	$69.08_{\pm0.01}$	$94.33_{\pm0.45}$	$82.54_{\pm 1.20}$	$43.13_{\pm0.71}$	$0.59_{\pm 0.15}$
	ViT	$58.77_{\pm 35.48}$	$71.47_{\pm 32.31}$	$63.97_{\pm 40.01}$	$48.68_{\pm 39.01}$	$5.63_{\pm 1.28}$	72.15 $_{\pm 0.00}$	$93.08_{\pm0.00}$	$78.26_{\pm0.00}$	$31.23_{\pm 0.00}$	$2.41_{\pm 0.00}$
	LightGBM	$47.24_{\pm 27.33}$	$78.04_{\pm 20.83}$	$63.45_{\pm 35.80}$	$64.10_{\pm 22.97}$	$1.30_{\pm 0.95}$	68.20 <sub>±4.63</sub>	95.68 <sub>±0.81</sub>	$81.69_{\pm 1.59}$	$45.04_{\pm 5.72}$	$0.61_{\pm 0.07}$
FAR	MLP	$47.59_{\pm 25.30}$	$84.04_{\pm 11.23}$	66.16+34 34	$64.40_{\pm 20.63}$	$1.14_{\pm 0.71}$	$63.92_{\pm 5.39}$	$87.10_{\pm 1.95}$	$64.22_{\pm 4.19}$	$46.45_{\pm 6.10}$	$1.40_{\pm 0.15}$
FAR	SVM	$41.86_{\pm 22.55}$	$79.07_{\pm 15.06}$	$62.27_{\pm 34.09}$	$69.93_{\pm 16.85}$	$1.27_{\pm 0.76}$	$66.18_{\pm 6.21}$	$93.44_{\pm 0.52}$	$75.46_{\pm 2.58}$	$44.26_{\pm 8.35}$	$1.23_{\pm 0.14}$
	XGBoost	$42.75 \pm 25.86$	$76.85_{\pm 16.49}$	$60.33_{\pm 35.88}$	$68.11_{\pm 20.43}$	$1.69_{\pm 0.57}$	54.88 <sub>±9.26</sub>	$83.44_{\pm 4.79}$	$65.03_{\pm 6.67}$	$56.68_{\pm 8.92}$	$1.41_{\pm 0.35}$
	DetectBERT	$45.01_{\pm 26.84}$	$76.37_{\pm 23.07}$	$60.56_{\pm 36.63}$	$64.57_{\pm 20.54}$	$3.15_{\pm 1.09}$	$73.23_{\pm 3.68}$	$96.74 \pm 0.17$	$83.15_{\pm 1.44}$	$35.67_{\pm 5.55}$	$0.99_{\pm 0.15}$
	ViT	$47.03_{\pm 29.55}$	$78.53_{\pm 20.17}$	$58.77_{\pm 37.82}$	$60.23_{\pm 22.66}$	$4.77_{\pm 1.20}$	$68.47_{\pm 3.94}$	$95.28_{\pm0.19}$	$73.69_{\pm 3.83}$	$27.71_{\pm 6.81}$	$3.81_{\pm 0.14}$

**Results.** Table 1 summarizes the performance of malware detectors on both LAMDA and API-Graph under the IID, NEAR, and FAR evaluation splits. All detectors perform strongly under IID conditions, but their effectiveness declines sharply as the temporal gap from training increases. For instance, LightGBM's F1-score on LAMDA drops from 97.49% (IID) to 59.48% (NEAR) and 47.24% (FAR), alongside a significant rise in the false negative rate, from 1.47% to 50.51% and 64.10%, respectively,—demonstrating increased difficulty. In contrast, the false positive rate (FPR) remains low and stable, likely due to the more consistent behavior of benign apps over time. Figure 2(a) further visualizes this trend, showing how F1-scores decline over time. Notably, we observe a sharp drop in performance between 2016 and 2017, indicating a significant distributional shift. A similar decline is evident from 2023 to 2024. In contrast, F1-scores increase from 2018 to the 2019–2022 period, suggesting that these intermediate years exhibit less drift relative to 2017 and 2018.

In APIGraph, LightGBM's F1-score drops from 85.95% (IID) to 66.77% (NEAR), but stabilizes at 68.20% on FAR. The F1-scores over the years in Figure 2(b) indicate a smaller degree of temporal drift, with only modest changes between years. Compared to the APIGraph, LAMDA shows a higher standard deviation in both NEAR and FAR, suggesting more pronounced and variable distributional shifts. Complex transformer-based models like DetectBERT and ViT also degrade with increasing drift severity. This supports our claim that LAMDA introduces stronger concept drift, making it a more challenging and realistic benchmark for evaluating long-term malware detection.

#### 4.2 VISUAL ANALYSIS OF CONCEPT DRIFT

**Setting.** To track how malware and benign class distributions evolve over time, we use two complementary visualization techniques: Jeffreys divergence heatmaps and t-SNE projections. Jeffreys divergence (Jeffreys, 1946), a symmetric information-theoretic measure, quantifies distributional shifts in individual static features across years. We compute pairwise divergences for all yearly combinations in LAMDA (2013–2025) and APIGraph (2012–2018). In parallel, we apply t-SNE (Van der Maaten & Hinton, 2008) to project high-dimensional feature vectors into 2D space. For direct comparison, we focus on four years shared by both datasets (2013, 2014, 2016, 2017), following common practice in prior malware drift studies (Pendlebury et al., 2019). Full year-wise t-SNE visualizations are provided in Appendix B.

**Analysis.** Figure 3 shows Jeffreys divergence heatmaps for both datasets. In both LAMDA and APIGraph, divergence increases as the temporal gap widens, confirming non-trivial concept drift. LAMDA exhibits a broader divergence range, particularly from 2022–2025, reflecting substantial feature distribution changes likely driven by evolving APIs, shifting development practices, and emerging malware behaviors. In contrast, APIGraph remains relatively stable, with limited divergence in its later years. Figure 4 provides t-SNE projections for the selected years. In LAMDA, samples become more dispersed by 2016–2017, suggesting increasing structural sparsity in feature space. Since t-SNE distorts global distances, we corroborate these patterns with Jeffreys divergence, which confirms genuine distributional shifts. APIGraph, by comparison, maintains tight and relatively static clusters, indicating limited structural evolution. Together, these trends underscore LAMDA's value as a temporally rich benchmark for studying real-world concept drift.

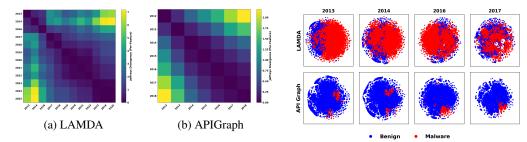


Figure 3: Jeffreys divergence heatmaps across Figure 4: t-SNE projections showing feature years for LAMDA and APIGraph.

space evolution for LAMDA and APIGraph.

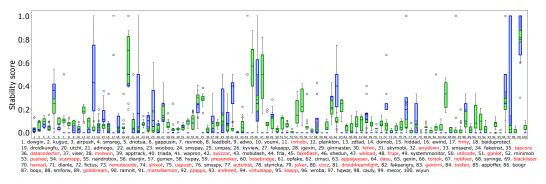


Figure 5: The distribution of feature stability scores for top 100 malware families. 58 families are common in both LAMDA (green) and APIGraph (blue) datasets, and families marked as red labels along x-axis available in LAMDA with minimum family size criteria.

#### COMPREHENSIVE DRIFT ANALYSIS

In this section, we present a more in-depth analysis of different types of drift, including feature stability (Section 5.1), temporal drift (Section 5.2), SHAP-based explainability (Section 5.3), and label drift (Section 5.4)

# FEATURE SPACE STABILITY ANALYSIS ON TOP MALWARE FAMILIES

**Analysis Setting.** We evaluate the temporal consistency of malware families using two complementary metrics: stability scores and Optimal Transport Dataset Distance (OTDD) (Alvarez-Melis & Fusi, 2020), following prior work (Zhang et al., 2020; Civitarese et al., 2022). The analysis is conducted in the original feature space and focuses on the 100 malware families with the largest sample sizes. Within each family, samples are temporally ordered from 2013 to 2025, though not all families contain data for every year. Following (Zhang et al., 2020), we partition each family's samples into ten equal subsets, each representing 10% of the family's total. For APIGraph, we identify 58 families with at least 10 samples between 2013 and 2018, meeting this subdivision requirement. Stability scores are then computed using Jaccard similarity across the ten subsets for both LAMDA and APIGraph.

**Stability Scores Analysis.** Figure 5 shows the distribution of consecutive pairwise stability scores across ten groups for each of the top 100 malware families. The number of samples per family varies considerably, ranging from 186 to 32,475, with a mean of 1,984 and a median of 535. The green box plots correspond to LAMDA, while the blue box plots represent APIGraph. Both datasets capture the temporal evolution of malware families, as reflected in the spread and median of stability scores. Broader spreads and lower medians in both datasets indicate greater behavioral variability over time. Notably, LAMDA includes more families and reflects broader evolutionary patterns than APIGraph. These differences suggest that detection models trained on LAMDA may offer improved insight into concept drift, benefiting from greater sample diversity and family coverage.

Table 2: This table shows label drift for Android malware samples, highlighting shifts in detection over time. **TS**: Total # of Malware Samples,  $B_C$ : Currently Labeled as Benign,  $\%B_C$ : Percentage of Total Malware Samples Currently Labeled as Benign.  $D_{Improved}$ : Improved Detection,  $D_{Weakened}$ : Weakened Detection,  $D_{Unchanged}$ : Unchanged Detection,  $DS_{Drop}$ : Significant Drop of Detection Count,  $DS_{Improve}$ : Detection Count Significantly Increased.

Year	TS	$B_C$	$ \%B_C $	$D_{Improved}$	$D_{Weakened}$	$D_{Unchanged}$	$DS_{Drop}$	$DS_{Increase}$
2013	44383	24	0.05	40436	439	3484	85	34481
2014	45756	345	0.75	37108	1554	6749	863	27239
2016	45134	177	0.39	26963	7485	10509	1160	13581
2017	21359	1108	5.19	7765	10289	2197	5061	3362
2018	39350	1242	3.16	17561	15346	5201	7304	7600
2019	41585	22	0.05	22905	9294	9364	467	7518
2020	46355	25	0.05	20755	3931	21644	294	8001
2021	35627	23	0.06	10385	4482	20737	176	2531
2022	41648	4	0.01	10445	3629	27570	121	2719
2023	7892	15	0.19	1763	1979	4135	592	416
2024	794	0	0.00	74	319	401	79	19
2025	23	0	0.00	6	2	15	0	2

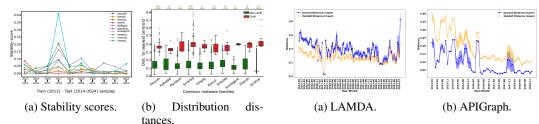


Figure 6: Stability and distribution analysis on malware families.

Figure 7: SHAP-based explanation drift on LAMDA and APIGraph datasets.

# 5.2 TEMPORAL DRIFT ANALYSIS ON COMMON MALWARE FAMILIES

**Analysis Setting.** We assess the drifting behavior over the years for the common families present from 2013 to 2025. We observe that *only 10* families appear consistently each year, except for 2025. Subsequently, we compute the year-wise stability score for the original feature set within each of these 10 common families. Additionally, we measure the distribution distances based on the CADE (Yang et al., 2021a) latent features in the test set. This experiment uses 2013 dataset for training and 2014 to 2024 samples serve as test sets.

Feature-Based Stability Evaluation. Figure 6a shows Jaccard similarity-based stability scores across consecutive yearly sample sets for 10 common malware families. Flatter curves indicate stronger temporal consistency, while sharp variations reflect feature drift. Most families such as airpush, dianjin, and smsreg exhibit relatively stable trends, suggesting consistent feature distributions over time. In contrast, families like umpay, fakeapp, and dnotua display pronounced fluctuations, with a notable spike around 2017–2018. The especially large peak for umpay signals substantial temporal drift, likely tied to evolving malware behaviors during that period. Overall, these results show that while many families retain stable characteristics, others undergo significant shifts, underscoring the importance of dynamic adaptation in detection models.

Latent Space Drift Detection via Distance Metrics. Figure 6b shows the distribution of distances from test samples (2014–2024) to their nearest class centroids across 10 common families, computed in the contrastive latent space (Yang et al., 2021d). Each test sample is encoded using the trained contrastive autoencoder, and Euclidean distances to class centroids are calculated. These distances are then normalized within each class using the *Median Absolute Deviation* (MAD). A sample is classified as *drifted* if its normalized MAD score  $A^{(k)}$  exceeds the empirical threshold  $T_{\rm MAD} = 3.5$ ; otherwise, it is considered *non-drifted*. This criterion identifies samples that deviate significantly from learned class distributions as potential drift instances. The resulting boxplots reveal a clear

Table 3: Comparison of concept drift adaptation methods on LAMDA and APIGraph under different labeling budgets. Results are reported as mean $_{\pm std}$ .

Budget	Method		APIGraph			LAMDA	
Dauger		F1 (%)	FNR (%)	FPR (%)	F1 (%)	FNR (%)	FPR (%)
50	Chen et al. (2023) CADE Yang et al. (2021d) TRANSCENDENT Barbero et al. (2022b)	$89.26_{\pm 0.31}$ $83.92_{\pm 1.75}$ $38.42_{\pm 1.52}$	$15.13_{\pm 0.31}$ $19.79_{\pm 2.09}$ $69.47_{\pm 1.83}$	$0.51_{\pm 0.03}$ $1.07_{\pm 0.15}$ $1.65_{\pm 0.27}$	$37.43_{\pm 2.43}$ $34.20_{\pm 1.80}$ $32.00_{\pm 1.50}$	$3.54_{\pm 0.32}$ $44.00_{\pm 2.50}$ $37.80_{\pm 2.10}$	$93.79_{\pm 2.71}$ $56.80_{\pm 2.20}$ $63.70_{\pm 1.80}$
100	Chen et al. (2023) CADE Yang et al. (2021d) TRANSCENDENT Barbero et al. (2022b)	$\begin{array}{c} \textbf{90.70}_{\pm 0.15} \\ 87.08_{\pm 0.54} \\ 40.17_{\pm 1.49} \end{array}$	$\begin{array}{c} \textbf{13.18}_{\pm 0.26} \\ 14.39_{\pm 1.08} \\ 68.03_{\pm 1.67} \end{array}$	$\begin{array}{c} \textbf{0.45}_{\pm 0.00} \\ 1.10_{\pm 0.16} \\ 1.42_{\pm 0.25} \end{array}$	$38.50_{\pm 2.20}$ $37.30_{\pm 1.60}$ $35.80_{\pm 1.30}$	$\begin{array}{c} \textbf{3.30}_{\pm 0.30} \\ 37.20_{\pm 2.10} \\ 28.40_{\pm 1.90} \end{array}$	$\begin{array}{c} 93.50_{\pm 2.60} \\ \textbf{45.20}_{\pm 2.00} \\ 70.10_{\pm 1.60} \end{array}$
200	Chen et al. (2023)   CADE Yang et al. (2021d)   TRANSCENDENT Barbero et al. (2022b)	$\begin{array}{c} \textbf{91.70}_{\pm 0.38} \\ 88.87_{\pm 0.24} \\ 42.38_{\pm 1.36} \end{array}$	$\begin{array}{c} \textbf{11.43}_{\pm 0.61} \\ 14.05_{\pm 0.35} \\ 66.48_{\pm 1.54} \end{array}$	$\begin{array}{c} \textbf{0.44}_{\pm 0.01} \\ 0.74_{\pm 0.01} \\ 1.15_{\pm 0.21} \end{array}$	$\begin{array}{c} \textbf{41.00}_{\pm 1.80} \\ 38.50_{\pm 1.30} \\ 39.00_{\pm 1.20} \end{array}$	$\begin{array}{c} \textbf{2.85}_{\pm 0.28} \\ 31.80_{\pm 1.80} \\ 21.30_{\pm 1.50} \end{array}$	$\begin{array}{c} 91.20_{\pm 2.30} \\ \textbf{38.50}_{\pm 1.60} \\ 76.30_{\pm 1.40} \end{array}$
400	Chen et al. (2023) CADE Yang et al. (2021d) TRANSCENDENT Barbero et al. (2022b)	$\begin{array}{c} 92.39_{\pm 0.29} \\ 89.16_{\pm 0.53} \\ 43.97_{\pm 1.21} \end{array}$	$\begin{array}{c} \textbf{10.18}_{\pm 0.52} \\ 13.23_{\pm 0.56} \\ 65.41_{\pm 1.42} \end{array}$	$\begin{array}{c} \textbf{0.45}_{\pm 0.02} \\ 0.80_{\pm 0.15} \\ 0.98_{\pm 0.17} \end{array}$	$43.00_{\pm 1.60}$ $45.40_{\pm 1.10}$ $40.60_{\pm 1.10}$	$2.50_{\pm 0.25}$ $59.20_{\pm 1.50}$ $15.60_{\pm 1.30}$	$89.80_{\pm 2.10}$ $10.10_{\pm 1.20}$ $82.80_{\pm 1.20}$

separation between non-drifted (green) and drifted (red) samples across families. Drifted samples consistently exhibit higher distances, with pronounced separation for families such as plankton, umpay, dianjin, and dnotua. In contrast, non-drifted samples cluster tightly around their centroids, indicating intra-family stability. Additional analyses are provided in Appendix C.

#### 5.3 TEMPORAL ANALYSIS OF SHAP-BASED EXPLANATION DRIFT

Analysis Setting. Explanation drift happens when the features a malware detector relies on for decisions change over time, even if its accuracy looks stable. To study this in the LAMDA and API-Graph (Zhang et al., 2020) datasets, we measure two types of change using SHAP (Lundberg & Lee, 2017) feature attributions: Jaccard distance, which checks how much the set of important features overlaps across time, and Kendall distance, which checks how consistent their rankings remain. Small distances mean the model reasons consistently, while large Jaccard or low Kendall values indicate shifts in reasoning. We compute SHAP values using KernelExplainer (Lundberg & Lee, 2017) with 100 background and 100 test samples per month. We calculate distances over the top-1000 features. To make figures easier to read, the x-axis in Figure 7a is labeled every three months, covering June 2013–January 2025 (shown from September 2013–December 2024). For APIGraph, the period is January 2013–December 2018 with the same labeling scheme (Figure 7b).

Jaccard and Kendall Distance Analysis. Figure 7a shows that LAMDA exhibits consistently high Jaccard distances (around 0.9), indicating substantial variability in the features the model relies on over time. A sharp dip around September 2017 marks a rare period of stability or possibly an anomaly. The Kendall distances reveal a moderate but steady pattern, reinforcing that both the set and the order of important features fluctuate significantly across time. In contrast, APIGraph (Figure 7b) shows a gradual downward trend in both measures, suggesting that its feature importance remains relatively stable. Overall, SHAP-based analysis highlights LAMDA's volatility compared to APIGraph's stability, underscoring LAMDA's value for studying concept drift, continual learning, and model robustness in dynamic malware detection scenarios (Zhang et al., 2020).

#### 5.4 LABEL DRIFT ANALYSIS ACROSS YEARS

Analysis Setting. We study how malware sample labels change over time, focusing on cases where a sample initially classified as malicious based on VirusTotal (VT) consensus is later reclassified as benign in a subsequent year, and vice versa. To analyze label drift, we use metadata from both AndroZoo (AZ) and VirusTotal (VT). AndroZoo provides metadata indicating how many VT engines flagged an application as malware at a given point. We first collected this metadata for a year. Then, for the same set of samples, we retrieved updated reports directly from VT in a later year. By comparing the two reports, we track how many samples changed their labels over time.

**Results.** Table 2 summarizes how labels change over time according to VT and AZ. The table reports yearly statistics from 2013 to 2025, including the total number of malware samples (TS) contained in LAMDA for each year. Among these, we report: the number of samples whose labels

have changed to benign  $(B_C)$ ; the samples where more VT engines now flag them as malware  $(D_{Improved})$ ; the samples where fewer VT engines now flag them as malware  $(D_{Weakened})$ ; and the samples where the total number of VT detections remains unchanged  $(D_{Unchanged})$ . For example, in 2013, among 44,383 malware samples that were initially detected as malicious, 24 are no longer flagged by any VT engine. Of the remaining 40,436 samples, 439 are now flagged by fewer engines, 3,484 show no change, and the rest are flagged by more engines. In addition, some samples show a significant drop  $(DS_{Dropped})$  or increase  $(DS_{Increase})$  in detection count, with these columns specifically highlighting drastic change (greater than 50%) in detection count.

## 6 CONCEPT DRIFT ADAPTATION

 In this section, we evaluate concept drift adaptation (CDA) using three state-of-the-art (SOTA) techniques: Chen et al. (2023), CADE (Yang et al., 2021d), and TRANSCENDENT (Barbero et al., 2022b), on the LAMDA dataset to highlight the unique adaptation challenges it poses compared to APIGraph. We adopt the Chen et al. (2023) active learning framework, which operates on a monthly cycle with labeling budgets of 50, 100, 200, and 400 samples. In each cycle, a subset of test samples is selected for labeling by human analysts, after which the model is retrained to mitigate performance degradation.

**Results.** Table 3 summarizes the performance of the evaluated methods. Across most labeling budgets, Chen-AL delivers the strongest results, except at the budget of 400 on LAMDA. Compared to APIGraph, LAMDA introduces substantially more *challenging adaptation scenarios*. While Chen et al. (2023) consistently outperforms CADE and TRANSCENDENT, all methods fail to generalize effectively on LAMDA. For instance, with a labeling budget of 400, Chen et al. (2023) achieves an F1-score of about 92% on APIGraph, but drops sharply to 43% on LAMDA. Similarly, CADE and TRANSCENDENT achieve only 45% and 40% F1-score, respectively. These findings empirically demonstrate that existing SOTA CDA techniques are insufficient when faced with longitudinal and complex distributional shifts such as those captured by LAMDA, highlighting the need for more advanced CDA approaches capable of adapting to long-term drift.

# 7 DISCUSSION AND LIMITATION

LAMDA's temporal structure facilitates evaluation of generalization under distributional shift, while its family diversity allows for studying malware evolution and model adaptation with limited data, including transfer and few-shot learning. Beyond these uses, LAMDA is also designed to facilitate continual learning research for both Domain and Class incremental learning (Rahman et al., 2022; Park et al., 2025), with detailed experimental results provided in Appendix I.

**Limitations.** While LAMDA provides a solid foundation for studying concept drift in malware analysis, we acknowledge a few limitations. First, it relies exclusively on DREBIN (Arp et al., 2014) features. Other static features such as control flow graphs and clustered API calls are out of scope of this work. Furthermore, we do not perform feature extraction on dynamic behaviors of the APKs which are observable only at runtime. Second, while previous work suggests a 10:90 malware-to-benign ratio, LAMDA attempts to maintain a 50:50 ratio, which may be viewed as downplaying the role of benign software distributions. However, this design emphasizes family diversity and balanced classes to create a challenging benchmark.

#### 8 CONCLUSION

In this paper, we present LAMDA, a large-scale, temporally structured Android malware dataset spanning over a decade. It enables analysis of how detection performance evolves with shifts in malware behavior and feature distributions. Our evaluations on supervised learning, feature stability, and explanation analysis demonstrate the impact of these temporal shifts. With broad coverage, diverse families, and static features, LAMDA offers a reproducible benchmark for advancing resilient and adaptive malware detection systems.

# ETHICS STATEMENT

This work does not involve human subjects or personally identifiable information. All Android applications were collected from the publicly available AndroZoo repository and contain no user data. To avoid misuse, we provide only extracted feature representations and associated metadata. As we didn't share live malware, the dataset can be made openly available to the research community. We believe the release of LAMDA fully complies with the ICLR Code of Ethics.

#### REPRODUCIBILITY STATEMENT

We provide comprehensive documentation and scripts to facilitate reproducibility. The dataset is publicly available on Zenodo in both Parquet and NPZ formats, accompanied by detailed metadata, allowing researchers to select the format most suitable for their needs. All code and experimental pipelines used in this work are released through an anonymous GitHub repository, which includes clear execution scripts and step-by-step guidelines to reproduce the reported results. Importantly, our source code is designed to be flexible, enabling researchers to incorporate new samples into the dataset with minimal effort.

In addition, we include detailed documentation of the dataset construction process, metadata, preprocessing steps, and experimental configurations. Together, these resources ensure that the findings of this paper can be independently verified, reproduced, and extended by the research community.

## REFERENCES

- Marco Alecci, Pedro Jesús Ruiz Jiménez, Kevin Allix, Tegawendé F Bissyandé, and Jacques Klein. AndroZoo: A Retrospective with a Glimpse into the Future. In *International Conference on Mining Software Repositories (MSR)*, 2024.
- Kevin Allix, Tegawendé F Bissyandé, Jacques Klein, and Yves Le Traon. AndroZoo: Collecting Millions of Android Apps for the Research Community. In *International Conference on Mining Software Repositories (MSR)*, 2016.
- David Alvarez-Melis and Nicolo Fusi. Geometric dataset distances via optimal transport. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Hyrum S Anderson and Phil Roth. Ember: An open dataset for training static pe malware machine learning models. In *Proceedings of the AAAI Workshop on Artificial Intelligence for Cyber Security (AICS)*, 2018a. URL https://github.com/elastic/ember. Available at https://github.com/elastic/ember.
- Hyrum S Anderson and Phil Roth. EMBER: An open dataset for training static PE malware machine learning models. *arXiv:1804.04637*, 2018b.
- 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 *Network and Distributed System Security Symposium (NDSS)*, 2014.
- AV-TEST. Malware statistics and trends report. https://www.av-test.org/en/statistics/malware/, 2025.
- Federico Barbero, Feargus Pendlebury, Fabio Pierazzi, and Lorenzo Cavallaro. Transcending transcend: Revisiting malware classification in the presence of concept drift. In *IEEE Symposium on Security and Privacy (S&P)*, 2022a.
- Federico Barbero, Feargus Pendlebury, Fabio Pierazzi, and Lorenzo Cavallaro. Transcending Transcend: Revisiting malware classification in the presence of concept drift. In *IEEE Symposium on Security and Privacy (S&P)*, 2022b.
- Brut. Apktool. https://apktool.org/, 2025. Accessed: 2025-04-26.
- Sen Chen, Minhui Xue, Zhushou Tang, Lihua Xu, and Haojin Zhu. Stormdroid: A streaminglized machine learning-based system for detecting android malware. In *ACM ASIA Conference on Computer and Communications Security (AsiaCCS)*, 2016.

- Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. *Proceedings of the* 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp.
   785–794, 2016.
  - Yizheng Chen, Zhoujie Ding, and David Wagner. Continuous learning for android malware detection. In *USENIX Security Symposium*, 2023.
    - Theo Chow, Zeliang Kan, Lorenz Linhardt, Lorenzo Cavallaro, Daniel Arp, and Fabio Pierazzi. Drift forensics of malware classifiers. In *ACM Workshop on Artificial Intelligence and Security (AISec)*, 2023.
    - Andrea Civitarese, Luca Demetrio, Fabio Pierazzi, and Lorenzo Cavallaro. Anoshift: A distribution shift benchmark for unsupervised anomaly detection under label scarcity. In *Advances in Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track*, 2022.
    - Dragos Georgian Corlatescu, Alexandru Dinu, Mihaela Gaman, and Paul Sumedrea. Embersim: A large-scale databank for boosting similarity search in malware analysis. In *Advances in Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track*, 2023.
    - Marius Dragoi, Elena Burceanu, Emanuela Haller, Andrei Manolache, and Florin Brad. Anoshift: A distribution shift benchmark for unsupervised anomaly detection. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
    - Scott Freitas, Yuxiao Dong, Joshua Neil, and Duen Horng Chau. Malnet: A large-scale malware network dataset for graph representation learning. In *Advances in Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track*, 2020.
    - Daniele Ghiani, Daniele Angioni, Angelo Sotgiu, Maura Pintor, and Battista Biggio. Understanding regression in continual learning for malware detection. 2025.
    - Andy Greenberg. Android ransomware's evolution is worrying researchers, 2020. URL https://www.wired.com/story/android-ransomware-worrying-evolution. Accessed: 2025-05-06.
    - Richard Harang and Ethan M. Rudd. Sorel-20m: A large scale benchmark dataset for malicious pe detection, 2020.
    - Harold Jeffreys. An invariant form for the prior probability in estimation problems. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, 186(1007):453–461, 1946.
    - Kaspersky. Mobile threat report 2024. https://securelist.com/mobile-threat-report-2024/115494/, 2024. Accessed: 2025-04-18.
    - Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, 30, 2017.
    - Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory G. Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2019.
    - Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
    - Enrico Mariconti, Lucky Onwuzurike, Panagiotis Andriotis, Emiliano De Cristofaro, Gordon Ross, and Gianluca Stringhini. MAMADROID: Detecting android malware by building markov chains of behavioral models. In *Network and Distributed System Security Symposium (NDSS)*, 2017.
    - The Hacker News. New android trojan 'soumnibot' evades detection by obfuscating manifest file, 2024. URL https://thehackernews.com/2024/04/new-android-trojan-soumnibot-evades.html. Accessed: 2025-05-06.

- Diane Oyen, Michal Kucer, Nicolas Hengartner, and Har Simrat Singh. Robustness to label noise depends on the shape of the noise distribution. *Advances in Neural Information Processing Systems*, 35:35645–35656, 2022.
  - Jimin Park, AHyun Ji, Minji Park, Mohammad Saidur Rahman, and Se Eun Oh. MalCL: Leveraging gan-based generative replay to combat catastrophic forgetting in malware classification. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2025.
  - Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12: 2825–2830, 2011.
  - Feargus Pendlebury, Fabio Pierazzi, Roberto Jordaney, Johannes Kinder, and Lorenzo Cavallaro. TESSERACT: Eliminating experimental bias in malware classification across space and time. In *USENIX Security Symposium*, 2019.
  - Mohammad Saidur Rahman, Scott E. Coull, and Matthew Wright. On the limitations of continual learning for malware classification. In *First Conference on Lifelong Learning Agents (CoLLAs)*, 2022.
  - Mohammad Saidur Rahman, Scott Coull, Qi Yu, and Matthew Wright. MADAR: Efficient continual learning for malware analysis with diversity-aware replay. *arXiv preprint arXiv:2502.05760*, 2025.
  - David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
  - Manuel Sebastián, Raul Rivera, Paraskevas Kotzias, and Juan Caballero. Avclass: A tool for massive malware labeling. In *International Symposium on Research in Attacks, Intrusions and Defenses (RAID)*, 2016.
  - Ravid Shwartz-Ziv, Micah Goldblum, Yucen Li, C Bayan Bruss, and Andrew G Wilson. Simplifying neural network training under class imbalance. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
  - Gido M van de Ven, Tinne Tuytelaars, and Andreas S Tolias. Three types of incremental learning. *Nature Machine Intelligence*, 2022.
  - Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research (JMLR)*, 2008.
  - VirusTotal. VirusTotal Stats. https://www.virustotal.com/qui/stats, 2025.
  - Yunyang Xiong, Zhanpeng Zeng, Rudrasis Chakraborty, Mingxing Tan, Glenn Fung, Yin Li, Vikas Singh, Richard Socher, and Caiming Xiong. Nyströmformer: A nyström-based algorithm for approximating self-attention. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 35, pp. 14138–14148, 2021.
  - Ke Xu, Yingjiu Li, Robert Deng, Kai Chen, and Jiayun Xu. Droidevolver: Self-evolving android malware detection system. In *IEEE European Symposium on Security and Privacy (EuroS&P)*, 2019.
  - Limin Yang, Yizheng Chen, and Gang Wang. Cade: Detecting and explaining concept drift samples for security applications. In *USENIX Security Symposium*, 2021a.
  - Limin Yang, Aravind Ciptadi, Ilya Laziuk, Amin Ahmadzadeh, and Gang Wang. Bodmas: An open dataset for learning based temporal analysis of pe malware. In *IEEE Security and Privacy Workshops (SPW)*, 2021b.

- Limin Yang, Arridhana Ciptadi, Ihar Laziuk, Ali Ahmadzadeh, and Gang Wang. BODMAS: An open dataset for learning based temporal analysis of PE malware. In *IEEE Security and Privacy Workshops (SPW)*, 2021c.
- 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*, 2021d.
- 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 *ACM Conference on Computer and Communications Security (CCS)*, 2020.
- Zhihan Zhou, Jiangchao Yao, Feng Hong, Ya Zhang, Bo Han, and Yanfeng Wang. Combating representation learning disparity with geometric harmonization. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

#### OVERVIEW OF APPENDIX

The following appendices provide further information:

- 1. A Dataset Statistics: Year-wise malware/benign counts and family distributions.
- 2. B Feature Description: Overview of all static features extracted from APKs.
- C Additional Analysis of Concept Drift: Visual analysis of concept drift and feature space stability analysis using OTDD
- D Model Architectures and Detail Results: Architectures and extended evaluation on LAMDA variants.
- 5. **E Behind the Scenes: Practical Challenges**: Technical and operational challenges during dataset construction.
- F Effect of Label Noise in Training Data: Impact of different VirusTotal thresholds on labeling.
- G More analysis on Label Drift Across Years: Year-wise analysis of evolving VirusTotal labels.
- 8. **H Scalability of LAMDA**: Instructions for extending LAMDA with new samples using our codebase.
- 9. I Continual Learning on LAMDA: Class- and domain-incremental learning benchmarks.
- 10. **J Computational Resources**: Hardware and runtime configuration for dataset generation.

# A DATASET STATISTICS

LAMDA benchmark is constructed from a total of 1,008,381 Android APKs, comprising 369,906 malware samples and 638,475 benign samples. Table 4 summarizes the yearly distribution of both malware and benign APKs. To mitigate class imbalance during training, our initial goal was to collect approximately 50,000 malware and 50,000 benign samples per year; this target could not be met in certain years. Specifically, we were unable to collect sufficient samples for the years 2017, 2021, 2023, 2024, and 2025 due to our labeling criterion—requiring a VirusTotal detection count of 4 or more for malware—and the limited availability of up-to-date samples in the AndroZoo repository (Allix et al., 2016; Alecci et al., 2024). Additional constraints, such as corrupted downloads and decompilation failures, further reduced the effective sample count in those years. Despite these limitations, LAMDA remains the largest Android malware dataset to date in terms of both total sample count and temporal coverage.

Beyond binary labels, LAMDA also includes family-level annotations for malware samples. As shown in Table 5, the dataset spans 1,380 distinct malware families, offering rich diversity for future analysis. Additionally, 150,604 samples are singletons, belonging to families that appear only once in the dataset, representing rare or unique variants. Moreover, 2,985 samples are marked as

Table 4: Year-wise distribution of total, malware, and benign samples.

Year	<b>Total Samples</b>	Malware Samples	Benign Samples
2013	86,431	44,383	42,048
2014	101,183	45,756	55,427
2016	109,193	45,134	64,059
2017	99,144	21,359	77,785
2018	104,292	39,350	64,942
2019	91,050	41,585	49,465
2020	102,073	46,355	55,718
2021	81,155	35,627	45,528
2022	86,416	41,648	44,768
2023	54,354	7,892	46,462
2024	48,427	794	47,633
2025	44,663	23	44,640
Total	1,008,381	369,906	638,475

Table 5: Year-wise breakdown of malware family distributions in LAMDA.

Year	New	Existing	Valid Family	#of Singleton	#of Unknown
2013	213	0	213	1550	24
2014	91	140	231	2482	345
2016	179	196	375	5861	177
2017	88	119	207	9063	1108
2018	153	220	373	20579	1242
2019	259	376	635	18916	22
2020	141	447	588	30644	25
2021	43	252	295	30020	23
2022	161	490	651	24927	4
2023	37	187	224	5922	15
2024	14	50	64	626	0
2025	1	7	8	14	0
Total	1,380			150,604	2,985

Table 6: Distribution of unknown malware samples by VirusTotal detection count.

VT Detection	4	5	6	7	8	9	10	11	12	13	14	15	18	19	Total
# of Unknown Sample	1643	664	226	153	133	65	68	15	3	4	5	4	1	1	2,985

"unknown", where AVClass2 is unable to confidently assign a family label. Table 6 reports the VirusTotal (VirusTotal, 2025) detection counts for these unknown-labeled samples, offering insight into their potential threat level even in the absence of a family tag.

This comprehensive summary, encompassing both class labels and family-level information, supports a wide range of research directions, including supervised detection, rare variant modeling, family classification, and concept drift analysis across diverse malware behaviors.

#### FEATURE DESCRIPTION В

Built upon static analysis of Android APKs, LAMDA incorporates a broad spectrum of executionfree features based on the features of Drebin (Arp et al., 2014). Table 7 summarizes the key categories of static features used in LAMDA (Arp et al., 2014). These include declared components (e.g., services, activities), permissions (requested and used), intent filters, restricted or suspicious API calls, and embedded network indicators such as hardcoded IPs and URLs.

Table 7: Static Features and Their Descriptions.

Feature	Description
Requested permissions	Permissions declared in the manifest (e.g., CAMERA, BLUETOOTH) indicating intended access to sensitive resources.
Declared activities and services	Registered components of the application, providing insight into its structural and behavioral composition.
Broadcast receivers	Components that handle specific system or custom intents (e.g., BOOT_COMPLETED), often linked to persistence or event-driven behavior.
Hardware components	Device capabilities required by the app (e.g., camera, Bluetooth), implying functional intent.
Intent filters	Define the types of intents components can respond to; critical for modeling potential entry points.
Used permissions	Permissions referenced in the small code, reflecting actual permission usage.
Restricted API calls	APIs that are protected by system permissions or grant access to sensitive resources.
Suspicious API calls	APIs heuristically associated with malicious or abnormal behavior.
Embedded IP addresses and URL domains	Hardcoded network endpoints that may indicate command-and-control (C&C) servers or tracking mechanisms.

Each APK is converted into a binary feature vector using a bag-of-tokens representation. Tokens are derived from the presence or absence of the static properties listed in Table 7. Since each application typically uses only a small fraction of the global feature space, the resulting vectors are sparse and high-dimensional. To address this, we apply different VarianceThreshold feature selection (Pedregosa et al., 2011), resulting in three dataset variants with different dimensionalities and sizes. Table 8 summarizes these variants. The Baseline variant uses a threshold of 0.001 (Rahman et al., 2025) and yields 4,561 binary features. Increasing the threshold to 0.01 results in a smaller, more compressed feature space with 925 features, while lowering it to 0.0001 expands the feature space to over 25,000 features.

#### **B.1 LAMDA VARIANTS**

For the LAMDA dataset variants, we apply different thresholds using the <code>VarianceThreshold</code> (varTh) feature selector. In the baseline configuration (varTh = 0.001), we retain 4,561 features with a total in-memory size of 222 MB. For a more relaxed threshold (varTh = 0.0001), we preserve 25,460 features, resulting in a memory size of 554 MB. Conversely, applying a stricter threshold (varTh = 0.01) yields 915 features with a reduced storage size of 138 MB. These information are summarized in Table 8.

Table 8: Summary of Dataset Variants by Variance Threshold.

Variant	Threshold	# Metadata	# Binary Features	Size
Baseline	0.001	5	4561	222MB
var_thresh_0.0001	0.0001	5	25460	554MB
var_thresh_0.01	0.01	5	925	138MB

## C ADDITIONAL ANALYSIS OF CONCEPT DRIFT

In this section, we provide in detail analysis of LAMDA visualization using t-SNE and feature space stability analysis using Optimal Transport Dataset Distance (OTDD) (Alvarez-Melis & Fusi, 2020).

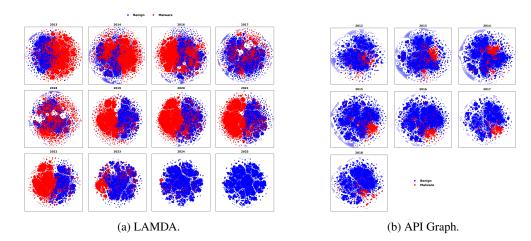


Figure 8: t-SNE projection of LAMDA and API Graph dataset, (a) t-SNE project of LAMDA from 2013 to 2025 (excluding 2015) and (b) t-SNE projection of API Graph from 2012 to 2018.

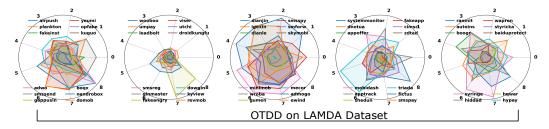


Figure 9: Optimal Transport Distance of 60 common families. Each of the plots shows the area of nine OTDD scores of 10 groups of 10 families in LAMDA.

#### C.1 VISUAL ANALYSIS OF CONCEPT DRIFT

To visualize the structural differences these features capture, we present t-SNE projections comparing LAMDA and API Graph (Dragoi et al., 2022) in Figure 8. LAMDA shows more scattered and diverse malware clusters over time, suggesting richer feature representations and stronger concept drift compared to the relatively compact structure in API Graph. This diversity, driven by the dynamic use of static tokens such as APIs and permissions, highlights the importance of broad and representative feature sets for modeling evolving malware behavior. Figure 13 further validates this hypothesis with varying number of virus total engine detection count.

#### C.2 FEATURE SPACE STABILITY ANALYSIS

**Optimal Transport Dataset Distance (OTDD) Analysis.** Figure 9 and 10 illustrates temporal distributional shifts using Optimal Transport Dataset Distance (OTDD) (Alvarez-Melis & Fusi, 2020), a geometric method for quantifying differences between probability distributions. To assess intra-family drift, we partition each malware family in the LAMDA and APIGraph datasets into ten chronological subsets and compute OTDD between consecutive pairs. The results are visualized via radar plots, where each axis represents a subset transition. Compact, regular shapes indicate temporal stability, while larger or irregular shapes signal drift. Comparing the two datasets, the LAMDA radar plots show both regular and irregular patterns indicating temporal shifts of malware families that causes concept drift. Similar behavior is also observed in the APIGraph dataset for the same families.

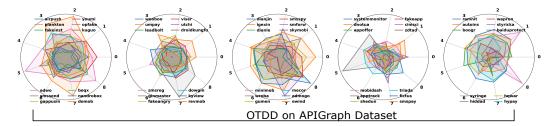


Figure 10: Optimal Transport Distance of 60 common families. Each of the plots shows the area of nine OTDD scores of 10 groups of 10 families in APIGraph.

# D ADDITIONAL EXPERIMENTAL DETAILS

In this section, we summarize the details of the model architectures and the training setup for each method used in our experiments. In addition, we present supplementary results.

#### D.1 DETAILS OF THE BASELINE METHODS

Multi-Layer Perceptron (MLP). The MLP model used for the experiments is adapted from prior work (Rahman et al., 2022; 2025) and is composed of four fully connected layers with the following sizes: 1024, 512, 256, and 128. Each hidden layer is followed by batch normalization, ReLU activation, and a dropout layer with a dropout rate of 0.5. The final output layer uses a sigmoid activation function for binary classification. The model is trained using Adam optimizer with a learning rate of 0.001, and batch size 512, as it stabilizes by this point, avoiding unnecessary GPU time.

**LightGBM.** In addition to MLP, we also use LightGBM (Ke et al., 2017), gradient-boosted decision tree ensemble, for binary classification. LightGBM is trained with up to 5000 estimators and a learning rate of 0.02, with early stopping based on Area Under the Curve (AUC) metric if no improvement is observed for 100 rounds. Each tree is allowed up to 256 leaves to provide high capacity for learning complex patterns. We apply 80% subsampling of both rows and features to mitigate overfitting. We also include  $L_1$  and  $L_2$  regularization to further penalize methods complexity to prevent overfitting. These hyperparameters are selected based on practices in malware detection benchmarks such as EMBER (Anderson & Roth, 2018a) and TESSERACT (Pendlebury et al., 2019).

**XGBoost.** The adapted XGBoost is configured with a tree depth of 12 and a learning rate of 0.05. We use log loss objective for binary classification (Chen & Guestrin, 2016; Anderson & Roth, 2018b). The method is trained for up to 3000 boosting rounds and uses the <code>gpu\_hist</code> tree construction method to accelerate training. The input data is loaded in XGBoost's <code>DMatrix</code> format, which is optimized for memory efficiency and fast training. We train the method on the full training data without applying early stopping and evaluate using log loss.

**Support Vector Machine (SVM).** A linear SVM model is implemented using LinearSVC and calibrated using CalibratedClassifierCV to enable probability outputs. This is essential for downstream evaluation where probabilistic thresholds or ranking-based metrics are used. Following prior work, the method is trained on the full dataset with a maximum of 10,000 iterations (Chen et al., 2023). Post-training, model memory usage is reported using psutil to assess resource footprint.

**detectBERT.** DetectBERT is a lightweight transformer-based model designed for malware classification that leverages DexBERT embeddings to learn full app-level representations. It projects the input vector into a hidden space using a fully connected layer, then prepends a learnable [CLS] token. The resulting sequence is processed by two stacked transformer layers utilizing Nyströmbased self-attention (Xiong et al., 2021), followed by LayerNorm and a classification head applied to the [CLS] output. The model also supports alternative aggregation strategies such as averaging or summation over token embeddings.

Table 9: Performance of models across IID, NEAR, and FAR splits for LAMDA on Baseline (VarianceThreshold = 0.001).

Split	Model	Accuracy	Precision	Recall	F1	ROC AUC	PR AUC	FPR	FNR
	LightGBM	97.74 ± 0.35	96.74 ± 0.31	98.26 ± 0.34	97.49 ± 0.17	99.55 ± 0.03	99.50 ± 0.11	$2.69 \pm 0.48$	1.74 ± 0.34
IID	MLP	$97.50 \pm 0.44$	$96.91 \pm 0.29$	$97.50 \pm 0.06$	$97.21 \pm 0.12$	$99.48 \pm 0.04$	$99.38 \pm 0.20$	$2.58 \pm 0.85$	$2.50 \pm 0.06$
ш	SVM	$95.61 \pm 0.61$	$94.78 \pm 1.41$	$95.18 \pm 0.76$	$94.98 \pm 1.07$	$98.89 \pm 0.28$	$98.75 \pm 0.46$	$4.09 \pm 0.55$	$4.82 \pm 0.76$
	XGBoost	$97.36 \pm 0.15$	$96.32 \pm 0.70$	$97.80 \pm 0.43$	$97.05 \pm 0.14$	$99.15 \pm 0.16$	97.68 ± 1.16	$2.96 \pm 0.01$	$2.20 \pm 0.43$
	detectBERT	$95.83 \pm 0.15$	$95.04\pm1.12$	$95.51\pm0.60$	$95.27 \pm 0.86$	$98.93 \pm 0.27$	$98.62 \pm 0.79$	$3.96\pm0.14$	$4.49\pm0.60$
	ViT	$95.57 \pm 0.80$	$94.14 \pm 2.42$	$95.83 \pm 0.73$	94.97±1.59	$98.91 \pm 0.37$	$98.61 \pm 0.85$	$4.63 \pm 0.81$	$4.17\pm0.73$
	LightGBM	85.83 ± 3.96	90.36 ± 5.21	49.49 ± 30.82	59.48 ± 28.20	74.05 ± 23.76	70.18 ± 27.10	1.85 ± 0.95	50.51 ± 30.82
NEAD	MLP	$83.90 \pm 3.75$	$78.12 \pm 13.98$	$48.05 \pm 30.42$	$56.57 \pm 28.41$	$82.71 \pm 11.19$	$67.94 \pm 24.59$	$3.98 \pm 1.19$	$51.95 \pm 30.42$
NEAR	SVM	$82.08 \pm 3.11$	$72.56 \pm 18.38$	$44.38 \pm 28.88$	$52.91 \pm 28.40$	$75.18 \pm 17.98$	$62.53 \pm 29.82$	$4.71 \pm 0.97$	$55.62 \pm 28.88$
	XGBoost	$84.59 \pm 3.75$	$86.18 \pm 9.02$	$46.16 \pm 30.94$	$55.84 \pm 29.73$	$77.75 \pm 16.85$	$68.14 \pm 26.55$	$2.14 \pm 0.86$	$53.84 \pm 30.94$
	detectBERT	$84.21 \pm 4.81$	$78.61 \pm 19.05$	$50.08 \pm 37.55$	$59.11 \pm 34.70$	$73.86\pm29.73$	$66.95\pm37.60$	$4.05\pm0.84$	$49.92\pm37.55$
	ViT	$83.66 \pm 5.36$	$73.73 \pm 22.28$	$51.32 \pm 39.01$	$58.77 \pm 35.48$	$71.47 \pm 32.31$	$63.97 \pm 40.01$	$5.63 \pm 1.28$	$48.68 \pm 39.01$
	LightGBM	83.94 ± 10.61	74.65 ± 34.66	35.90 ± 22.97	47.24 ± 27.33	78.04 ± 20.83	63.45 ± 35.80	$1.30 \pm 0.95$	64.10 ± 22.97
FAR	MLP	$83.45 \pm 10.74$	$76.12 \pm 33.39$	$35.60 \pm 20.63$	$47.59 \pm 25.30$	$84.04 \pm 11.23$	$66.16 \pm 34.34$	$1.14 \pm 0.71$	$64.40 \pm 20.63$
rak	SVM	$80.99 \pm 11.98$	$72.89 \pm 35.60$	$30.07 \pm 16.85$	$41.86 \pm 22.55$	$79.07 \pm 15.06$	$62.27 \pm 34.09$	$1.27 \pm 0.76$	69.93 ± 16.85
	XGBoost	$82.03 \pm 11.07$	$70.00 \pm 37.26$	$31.89 \pm 20.43$	$42.75 \pm 25.86$	$76.85 \pm 16.49$	$60.33 \pm 35.88$	$1.69 \pm 0.57$	$68.11 \pm 20.43$
	detectBERT	$81.79 \pm 11.09$	$66.49 \pm 38.11$	$35.43\pm20.54$	$45.01\pm26.84$	$76.37 \pm 23.07$	$60.56 \pm 36.63$	$3.15\pm1.09$	$64.57 \pm 20.54$
	ViT	$81.98 \pm 9.44$	$63.56 \pm 38.55$	$39.77 \pm 22.66$	$47.03 \pm 29.55$	$78.53 \pm 20.17$	$58.77 \pm 37.82$	$4.77 \pm 1.20$	$60.23 \pm 22.66$

ViT. We utilize a ViT-based model adapts the Vision Transformer (ViT) framework to malware classification using static feature vectors. Each input sample, represented as a flat feature vector, is projected into a hidden-dimensional token via a linear embedding layer. A learnable <code>[CLS]</code> token is optionally prepended to the sequence, and positional embeddings are added to all tokens. The resulting token sequence is passed through a stack of seven Transformer encoder blocks, each consisting of LayerNorm, multi-head self-attention, and a two-layer feed-forward network with GELU activation and residual connections. After encoding, the final hidden state of the <code>[CLS]</code> token (or the feature token if <code>[CLS]</code> is not used) is extracted and passed through a classifier composed of a LayerNorm and a linear output layer. The model outputs logits over malware families and supports fine-grained malware classification or detection.

All models are trained on three different LAMDA variants with VarianceThreshold (VarTh)  $\in \{0.01, 0.001, 0.0001\}$  where VarTh = 0.001 is the baseline. No task-specific tuning or dataset-specific hyperparameter adjustments are performed to ensure fair comparisons across splits and datasets.

# D.2 BASELINE PERFORMANCE

We compare LAMDA baseline with API Graph (Zhang et al., 2020) dataset and provide a comprehensive results on four methods discussed above using AnoShift-style (Dragoi et al., 2022) splits. A subset of Table 9 and Table 10 are explained in the main body of the paper. We present the results with more performance metrics.

We compare the LAMDA baseline with the API Graph (Zhang et al., 2020) dataset and present comprehensive results using four models under AnoShift-style (Dragoi et al., 2022) splits. While a subset of results is highlighted in Table 9 and Table 10 in the main paper, we report extended metrics here for completeness.

Across both NEAR and FAR splits, LAMDA consistently exhibits lower scores across all performance metrics compared to API Graph, and notably higher false negative rates (FNR). These trends clearly indicate that LAMDA captures a significantly higher degree of concept drift. Furthermore, the standard deviation across metrics is substantially higher in LAMDA, especially for drifted years, underscoring the dataset's temporal instability in detection performance—validating the presence of concept drift.

# D.3 LAMDA VARIANTS AND DRIFT SENSITIVITY

LAMDA offers flexibility to researchers for Android malware analysis by supporting different feature selection variants. In this section, we evaluate two additional variants of LAMDA. As shown in Table 11 and Table 12, we report detailed performance results for the four methods and configurations used in the primary analysis of concept drift with <code>VarianceThreshold</code> of 0.001. The variant using a threshold of 0.01 exhibits a relatively higher average F1-score across NEAR and FAR

Table 10: Performance of models on across IID, NEAR, and FAR splits for API Graph.

Split	Model	Accuracy	Precision	Recall	F1	ROC AUC	PR AUC	FPR	FNR
	LightGBM	97.02 ± 0.00	95.78 ± 0.00	$73.44 \pm 0.00$	83.14 ± 0.00	98.93 ± 0.00	94.92 ± 0.00	$0.36 \pm 0.00$	26.56 ± 0.00
IID	MLP	$97.35 \pm 0.20$	$94.84 \pm 1.49$	$77.79 \pm 2.85$	$85.43 \pm 1.35$	$94.62 \pm 0.92$	$89.33 \pm 1.75$	$0.47 \pm 0.16$	$22.21 \pm 2.85$
ш	SVM	$96.64 \pm 0.00$	$94.12 \pm 0.00$	$70.85 \pm 0.00$	$80.84 \pm 0.00$	$97.27 \pm 0.00$	$90.90 \pm 0.00$	$0.49 \pm 0.00$	$29.15 \pm 0.00$
	XGBoost	$96.30 \pm 0.00$	$91.57 \pm 0.00$	$69.37 \pm 0.00$	$78.94 \pm 0.00$	$95.93 \pm 0.00$	$89.08 \pm 0.00$	$0.71 \pm 0.00$	$30.63 \pm 0.00$
	detectBERT	$97.01 \pm 0.00$	$94.59 \pm 0.00$	$74.40 \pm 0.00$	$83.05 \pm 0.00$	$98.82 \pm 0.00$	$94.19 \pm 0.00$	$0.48 \pm 0.00$	$25.60 \pm 0.00$
	ViT	$97.49\pm0.00$	$91.85 \pm 0.00$	$82.18\pm0.00$	$86.64 \pm 0.00$	$98.65 \pm 0.00$	$93.73\pm0.00$	$0.81\pm0.00$	$17.82\pm0.00$
	LightGBM	94.84 ± 0.00	94.07 ± 0.00	51.33 ± 0.00	66.42 ± 0.00	96.28 ± 0.00	83.45 ± 0.00	$0.36 \pm 0.00$	48.67 ± 0.00
NEAR	MLP	$94.66 \pm 0.34$	$88.22 \pm 3.83$	$53.67 \pm 4.84$	$66.55 \pm 3.19$	$85.97 \pm 1.25$	$72.24 \pm 2.46$	$0.81 \pm 0.37$	$46.33 \pm 4.84$
NEAR	SVM	$94.01 \pm 0.00$	$81.70 \pm 0.00$	$51.18 \pm 0.00$	$62.93 \pm 0.00$	$90.29 \pm 0.00$	$70.84 \pm 0.00$	$1.26 \pm 0.00$	$48.82 \pm 0.00$
	XGBoost	$92.76 \pm 0.00$	$77.60 \pm 0.00$	$38.11 \pm 0.00$	$51.12 \pm 0.00$	$92.13 \pm 0.00$	$65.74 \pm 0.00$	$1.21 \pm 0.00$	$61.89 \pm 0.00$
	detectBERT	$95.24 \pm 0.29$	$90.54 \pm 1.04$	$56.87 \pm 0.71$	$69.08 \pm 0.00$	$94.33 \pm 0.45$	$82.54 \pm 1.20$	$0.59 \pm 0.15$	$43.13 \pm 0.71$
	ViT	$94.73\pm0.00$	$76.15\pm0.00$	$68.77 \pm 0.00$	$72.15\pm0.00$	$93.08 \pm 0.00$	$78.26 \pm 0.00$	$2.41\pm0.00$	$31.23\pm0.00$
	LightGBM	95.31 ± 0.61	87.48 ± 2.48	57.41 ± 6.57	69.07 ± 4.73	96.11 ± 0.40	80.94 ± 1.17	$0.84 \pm 0.24$	42.59 ± 6.57
FAR	MLP	$94.50 \pm 0.84$	$82.31 \pm 6.28$	$51.14 \pm 7.90$	$62.80 \pm 6.68$	$90.02 \pm 2.91$	$71.38 \pm 5.89$	$1.12 \pm 0.43$	$48.86 \pm 7.90$
FAK	SVM	$94.33 \pm 0.60$	$77.38 \pm 0.33$	$54.34 \pm 7.27$	$63.56 \pm 5.08$	$92.73 \pm 0.92$	$71.46 \pm 2.43$	$1.60 \pm 0.20$	$45.66 \pm 7.27$
	XGBoost	$93.70 \pm 0.35$	$75.99 \pm 2.62$	$46.04 \pm 1.92$	$57.29 \pm 1.45$	$87.78 \pm 1.90$	$66.49 \pm 4.80$	$1.49 \pm 0.28$	$53.96 \pm 1.92$
	detectBERT	$95.85 \pm 0.53$	$86.50 \pm 1.92$	$64.33 \pm 5.55$	$73.23 \pm 3.68$	$96.74 \pm 0.17$	$83.15 \pm 1.44$	$0.99 \pm 0.15$	$35.67 \pm 5.55$
	ViT	$94.01 \pm 0.71$	$65.50 \pm 1.51$	$72.29 \pm 6.81$	$68.47 \pm 3.94$	$95.28 \pm 0.19$	$73.69 \pm 3.83$	$3.81 \pm 0.14$	$27.71 \pm 6.81$

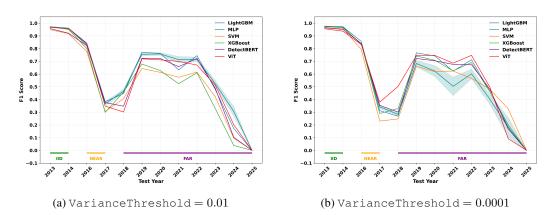


Figure 11: F1-scores on different models on based on AnoShift-style split on LAMDA. (a) for  $VarianceThreshold\ 0.01$  and (b) for  $VarianceThreshold\ 0.0001$ .

splits compared to that of the Baseline (0.001) and var\_thresh\_0.0001 variants. Figure 11 shows this trend as well, highlighting the comparative performance of the LAMDA variants.

To further understand how these feature selection variants influence drift sensitivity, we focus on the NEAR split with the SVM. For false positive rate (FPR), both the baseline and varTh=0.0001 maintain relatively low values ( $\sim 4.07 \pm 0.18$ ), whereas varTh=0.01 shows a sharp increase to  $17.09 \pm 2.87$ . This suggests that reducing the feature set too aggressively may misclassify benign as malware. Conversely, the false negative rate (FNR) slightly improves under varTh=0.01, decreasing from  $\sim 55.62 \pm 28.88$  and  $\sim 57.72 \pm 28.87$ , baseline and varTh=0.0001, respectively, to  $46.22 \pm 23.97$ . This indicates that even with fewer features, the model may still capture certain generalizable malware traits, improving detection of some malicious samples.

However, for the LightGBM model, this change is accompanied by a drop in precision. While in NEAR region both the baseline and varTh=0.0001 variants maintain high precision scores (around  $90.36 \pm 5.21$ ), the varTh=0.01 variant yields a lower precision of  $75.03 \pm 15.03$ . This reflects a shift in the methods decision behavior under more aggressive feature selection, emphasizing the importance of balancing dimensionality reduction with predictive consistency.

In summary, while varTh=0.01 may occasionally help with generalization under drift, it also amplifies misclassification of benign apps and reduces predictive stability. The baseline and varTh=0.0001 offer more shift in data distribution.

Table 11: Performance of models across IID, NEAR, and FAR splits for LAMDA variant of VarianceThreshold (0.01).

Split	Model	Accuracy	Precision	Recall	F1	ROC AUC	PR AUC	FPR	FNR
	LightGBM	97.08 ± 0.21	96.40 ± 0.22	96.95 ± 1.02	96.67 ± 0.61	99.42 ± 0.00	99.33 ± 0.13	2.94 ± 0.43	3.05 ± 1.02
IID	MLP	$96.89 \pm 0.32$	$96.11 \pm 0.37$	$96.84 \pm 1.26$	$96.47 \pm 0.67$	$99.29 \pm 0.08$	$99.15 \pm 0.18$	$3.18 \pm 0.54$	$3.16 \pm 1.26$
ш	SVM	$94.13 \pm 0.87$	$92.59 \pm 2.43$	$94.10 \pm 0.84$	$93.33 \pm 1.65$	$97.86 \pm 0.54$	$97.34 \pm 1.15$	$5.83 \pm 0.90$	$5.90 \pm 0.84$
	XGBoost	$96.65 \pm 0.60$	$95.59 \pm 1.37$	$96.74 \pm 0.80$	$96.16 \pm 1.09$	$99.01 \pm 0.26$	$98.00 \pm 0.36$	$3.45 \pm 0.43$	$3.26 \pm 0.80$
	detectBERT	$94.75 \pm 1.53$	$93.58 \pm 3.04$	$94.34 \pm 2.10$	$93.96 \pm 2.58$	$98.47 \pm 0.51$	$98.13 \pm 1.02$	$4.98 \pm 1.15$	$5.66 \pm 2.10$
	ViT	$96.66 \pm 0.07$	$96.02 \pm 0.37$	$96.42 \pm 0.93$	$96.22 \pm 0.65$	$99.07 \pm 0.38$	$98.80 \pm 0.82$	$3.24 \pm 0.59$	$3.58 \pm 0.93$
	LightGBM	84.17 ± 3.77	75.03 ± 15.03	54.00 ± 27.34	61.38 ± 24.22	73.58 ± 22.25	66.97 ± 27.37	5.86 ± 0.94	46.00 ± 27.34
NEAR	MLP	$83.39 \pm 3.52$	$73.64 \pm 16.21$	$52.23 \pm 26.78$	$59.77 \pm 24.33$	79.17 ± 14.79	$66.60 \pm 24.33$	$6.09 \pm 1.22$	$47.77 \pm 26.78$
NEAK	SVM	$75.73 \pm 6.47$	$54.52 \pm 25.76$	$53.78 \pm 23.97$	$54.13 \pm 24.86$	$74.88 \pm 14.32$	$58.72 \pm 25.70$	$17.09 \pm 2.87$	$46.22 \pm 23.97$
	XGBoost	$81.63 \pm 3.39$	$68.78 \pm 20.03$	$47.48 \pm 27.52$	$54.98 \pm 26.43$	$77.32 \pm 14.71$	$62.60 \pm 26.89$	$6.57 \pm 0.68$	$52.52 \pm 27.52$
	detectBERT	$82.23 \pm 5.19$	$67.95 \pm 24.83$	$54.66 \pm 34.89$	$59.95 \pm 31.61$	$71.17 \pm 30.38$	$63.83 \pm 35.68$	$8.79 \pm 0.97$	$45.34 \pm 34.89$
	ViT	$83.91 \pm 5.34$	$75.41 \pm 22.66$	$51.05 \pm 37.34$	$59.35 \pm 34.67$	$72.45 \pm 31.57$	$64.86 \pm 39.94$	$4.87 \pm 0.34$	$48.95 \pm 37.34$
	LightGBM	84.68 ± 9.54	71.15 ± 35.65	39.50 ± 23.10	50.13 ± 27.27	75.80 ± 21.64	61.66 ± 35.89	2.34 ± 1.72	60.50 ± 23.10
EAD	MLP	$85.15 \pm 9.23$	$76.93 \pm 33.45$	$41.38 \pm 21.41$	$53.00 \pm 25.53$	$87.00 \pm 10.48$	$68.52 \pm 33.91$	$1.80 \pm 1.61$	$58.62 \pm 21.41$
FAR	SVM	$77.50 \pm 10.05$	$55.00 \pm 32.97$	$36.40 \pm 16.42$	$42.16 \pm 23.50$	$76.19 \pm 12.22$	$52.48 \pm 31.86$	$9.28 \pm 3.56$	$63.60 \pm 16.42$
	XGBoost	$76.73 \pm 6.52$	$56.84 \pm 36.66$	$34.05 \pm 17.19$	$40.84 \pm 25.02$	$69.95 \pm 15.23$	$51.85 \pm 34.08$	$8.71 \pm 3.04$	65.95 ± 17.19
	detectBERT	$81.52 \pm 9.77$	$60.48 \pm 37.30$	$40.21 \pm 22.98$	$47.06 \pm 29.24$	$76.28 \pm 21.99$	$57.68 \pm 36.07$	$6.29 \pm 2.14$	$59.79 \pm 22.98$
	ViT	$82.84 \pm 10.71$	$67.79 \pm 38.23$	$38.47 \pm 21.47$	$47.74 \pm 27.67$	$77.55 \pm 24.71$	$61.82 \pm 37.79$	$2.85 \pm 1.43$	$61.53 \pm 21.47$

Table 12: Performance of models across IID, NEAR, and FAR splits for LAMDA variant of VarianceThreshold (0.0001).

Split	Model	Accuracy	Precision	Recall	F1	ROC AUC	PR AUC	FPR	FNR
	LightGBM	97.73 ± 0.04	96.80 ± 0.33	98.11 ± 0.12	97.45 ± 0.23	99.58 ± 0.00	99.54 ± 0.09	2.61 ± 0.25	1.89 ± 0.12
IID	MLP	$97.42 \pm 0.21$	$96.44 \pm 0.35$	$97.80 \pm 0.23$	$97.12 \pm 0.15$	$99.31 \pm 0.06$	$99.24 \pm 0.17$	$2.92 \pm 0.41$	$2.20 \pm 0.23$
ш	SVM	$96.68 \pm 0.02$	$96.68 \pm 0.08$	$95.77 \pm 0.72$	$96.22 \pm 0.40$	$99.13 \pm 0.20$	$99.13 \pm 0.29$	$2.69 \pm 0.48$	$4.23 \pm 0.72$
	XGBoost	$97.39 \pm 0.26$	$96.27 \pm 0.69$	$97.83 \pm 0.53$	$97.04 \pm 0.61$	$99.35 \pm 0.08$	$98.92 \pm 0.10$	$3.01 \pm 0.04$	$2.17 \pm 0.53$
	detectBERT	$96.24 \pm 0.23$	$95.33 \pm 0.90$	$96.15 \pm 0.87$	$95.74 \pm 0.89$	$98.92 \pm 0.40$	$98.61 \pm 0.93$	$3.76 \pm 0.28$	$3.85 \pm 0.87$
	ViT	$95.27 \pm 0.50$	$93.32 \pm 1.29$	$96.08 \pm 1.41$	$94.68 \pm 1.35$	$98.76 \pm 0.34$	$98.56 \pm 0.69$	$5.49 \pm 0.41$	$3.92 \pm 1.41$
	LightGBM	85.55 ± 3.91	90.24 ± 5.75	48.34 ± 30.74	58.46 ± 28.66	74.58 ± 23.25	70.47 ± 26.85	$1.70 \pm 0.80$	51.66 ± 30.74
NEAR	MLP	$84.22 \pm 3.38$	$80.23 \pm 13.34$	$49.38 \pm 27.83$	58.79 ± 25.81	81.65 ± 13.28	$69.02 \pm 25.24$	$3.71 \pm 0.95$	$50.62 \pm 27.83$
NEAK	SVM	81.79 ± 3.36	$71.61 \pm 21.55$	$42.28 \pm 28.87$	51.19 ± 29.64	$76.15 \pm 16.45$	$63.54 \pm 29.13$	$4.07 \pm 0.18$	57.72 ± 28.87
	XGBoost	84.58 ± 3.41	$86.82 \pm 6.81$	$46.86 \pm 30.72$	$56.40 \pm 28.84$	$77.29 \pm 18.06$	69.61 ± 25.51	$2.52 \pm 1.41$	$53.14 \pm 30.72$
	detectBERT	$84.30 \pm 4.99$	$77.92 \pm 19.62$	$50.98 \pm 37.46$	$59.73 \pm 34.36$	$74.34 \pm 28.90$	$67.51 \pm 36.87$	$4.33 \pm 0.78$	$49.02 \pm 37.46$
	ViT	$81.49 \pm 6.49$	$63.97 \pm 26.76$	$58.13 \pm 36.00$	$60.56 \pm 32.07$	$69.96 \pm 31.48$	$59.47 \pm 40.55$	$11.55 \pm 0.86$	$41.87 \pm 36.00$
	LightGBM	83.96 ± 10.66	75.34 ± 34.34	35.62 ± 23.11	47.02 ± 27.47	77.98 ± 21.70	64.21 ± 35.59	$1.18 \pm 0.85$	64.38 ± 23.11
FAR	MLP	$80.88 \pm 11.82$	$71.93 \pm 36.09$	29.64 ± 17.51	$40.99 \pm 23.11$	$82.52 \pm 12.80$	$63.53 \pm 35.27$	$1.47 \pm 1.11$	$70.36 \pm 17.51$
FAR	SVM	81.23 ± 11.94	$73.48 \pm 34.34$	$32.40 \pm 16.44$	$44.08 \pm 21.95$	$76.94 \pm 16.51$	$63.06 \pm 32.51$	$1.31 \pm 0.86$	$67.60 \pm 16.44$
	XGBoost	$83.47 \pm 10.41$	$73.53 \pm 35.83$	$35.43 \pm 20.83$	$46.97 \pm 25.75$	77.49 ± 19.29	$63.56 \pm 34.89$	$1.32 \pm 0.67$	$64.57 \pm 20.83$
	detectBERT	$82.54 \pm 10.53$	$67.37 \pm 38.73$	$37.34 \pm 21.46$	$46.62 \pm 27.85$	$77.47 \pm 22.32$	$61.65 \pm 37.23$	$2.92 \pm 1.21$	$62.66 \pm 21.46$
	ViT	$82.03 \pm 7.65$	$59.73 \pm 37.03$	$45.42 \pm 23.33$	$50.25 \pm 30.10$	$76.63 \pm 18.94$	$58.53 \pm 36.04$	$8.09 \pm 2.42$	$54.58 \pm 23.33$

## E BEHIND THE SCENES: PRACTICAL CHALLENGES IN LAMDA CREATION

#### E.1 ADMINISTRATIVE CHALLENGES

Downloading large volumes of real-world malware presents significant cybersecurity risks within any institutional environment. During our data collection process, the downloading and unpacking of live malware samples triggered internal threat detection systems, as automated security tools flagged these activities as potential breaches. To mitigate these risks, we implemented strict containment policies, such as disabling execution permissions. Additionally, we worked closely with the university's cybersecurity team to obtain the necessary approvals and ensure compliance with all relevant security policies. We maintained continuous communication with them throughout the process to ensure proper coordination and promptly address any emerging issues.

#### E.2 TECHNICAL CHALLENGES

We faced several technical constraints during the sample collection and processing pipeline. The AndroZoo platform imposes strict download rate limits, allowing only 40 concurrent downloads per user. As a result, we had to be extremely cautious to avoid violating their terms and conditions. Unfortunately, accidental oversights on our part led to temporary request blocks which disrupted the collection process. Similarly, the VirusTotal API has strict rate limits, which can significantly slow down the retrieval of metadata. Additionally, a considerable number of APKs failed to decompile successfully using Apktool. These failures were often due to obfuscation, corrupted files, or non-standard packaging formats. So, we had to perform multiple rounds of sampling to reach our target number of usable samples.

Table 13: Effect of thresholding on sample counts and relative percentage change (w.r.t. threshold 4).

Threshold	Benign Samples	Malware Samples	% Change (Benign)	% Change (Malware)
1	638,475	369,906	0.0%	0.0%
2	638,475	369,906	0.0%	0.0%
3	638,475	369,906	0.0%	0.0%
4	638,475	369,906	0.0%	0.0%
5	638,475	324,927	0.0%	↓ 12.16%
6	638,475	281,824	0.0%	↓ 23.81%
7	638,475	241,690	0.0%	↓ 34.66%
8	638,475	206,644	0.0%	↓ 44.14%
9	638,475	177,707	0.0%	↓ 51.96%
10	638,475	155,376	0.0%	↓ 58.00%
11	638,475	138,041	0.0%	↓ 62.68%
12	638,475	123,783	0.0%	↓ 66.53%
13	638,475	111,350	0.0%	↓ 69.89%

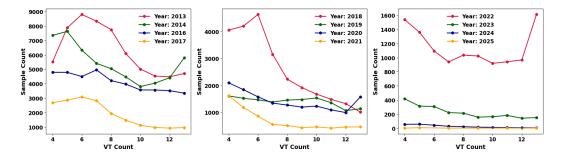


Figure 12: Sample Count Year-wise for Each VT Detection Threshold.

## F EFFECT OF LABEL NOISE IN TRAINING DATA

Label noise in Android malware family classification, particularly when using the Drebin (Arp et al., 2014) feature set can significantly impact model performance due to ambiguous and overlapping feature representations. As demonstrated in recent work (Oyen et al., 2022), the robustness of classification models depends not only on the amount of label noise but also on its distribution within the feature space. Specifically, feature-dependent label noise, where the probability of a label flip is contingent on the position of a sample in feature space, can cause a substantial drop in accuracy, even at low noise levels.

This is especially relevant for Drebin features, where different malware families may share static features like permissions ( $x_1 = \mathtt{INTERNET}, x_2 = \mathtt{SEND\_SMS}$ ), API calls ( $x_3 = \mathtt{getDeviceId}$ ), and hardware access ( $x_4 = \mathtt{READ\_PHONE\_STATE}$ ). Samples with minimal or ambiguous patterns (e.g.,  $x_5 = \mathtt{ACCESS\_NETWORK\_STATE}$  and  $x_6 = \mathtt{RECEIVE\_BOOT\_COMPLETED}$ ) are likely to fall near decision boundaries, increasing the risk of mislabeling. Such feature-dependent noise is more detrimental than uniform or class-dependent noise and warrants careful consideration in malware classification tasks.

Table 13 shows the effect of increasing the VirusTotal (VT) detection (VirusTotal, 2025) threshold on malware labeling in the LAMDA dataset. According to previous studies (Pendlebury et al., 2019; Rahman et al., 2022; Yang et al., 2021d), a sample is considered benign if  $vt\_detection = 0$ , and labeled as malware if  $vt\_detection \ge 4$ . As seen from the Benign sample column in the Table 13, the number of benign samples remains unchanged across all thresholds, since the benign definition is fixed and independent of the malware thresholding rule. However, the number of malware samples decreases significantly as the threshold increases from 4 to 13. For instance, at a threshold of 10, the number of malware samples drops by 58% compared to the baseline at threshold 4. This trend continues, reaching a 69.89% reduction at threshold 13. These results demonstrate that requiring stronger agreement among antivirus engines (i.e., a higher threshold) leads to more

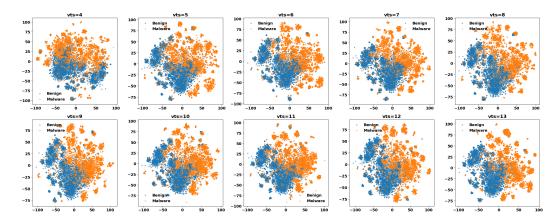


Figure 13: t-SNE visualization of benign and malware samples at varying Virus Total (VT) detection threshold.

conservative malware labeling, effectively excluding a substantial portion of potentially malicious samples. While this may improve the confidence in the labeled malware, it also drastically reduces dataset coverage. Therefore, the choice of VT threshold directly impacts the balance between label precision and data availability, and threshold 4 provides a practical trade-off commonly adopted in existing literature (Xu et al., 2019; Pendlebury et al., 2019; Park et al., 2025).

Figure 12 illustrates the distribution of malware sample counts across different VirusTotal (VT) threshold (VirusTotal, 2025) values for each year from 2013 to 2025 (except 2015). The VT count, plotted on the x-axis, represents the number of antivirus (AV) engines that flagged a sample as malicious, serving as a metric for detection consensus or confidence. A consistent trend is observed across all years — as the VT threshold increases, the number of flagged samples decreases. This suggests that only a small fraction of malware samples achieve strong consensus among AV engines, while the majority are detected by relatively few engines. The sample count is highest between VT counts of 5 to 7, especially in earlier years such as 2013–2017, indicating a moderate level of agreement in those periods. In contrast, from 2022 onward, the overall volume of detected samples decreases sharply, and the detections are largely concentrated in the lower VT ranges, which may reflect advancements in malware evasion techniques or shifts in detection criteria. These observations justify the use of a VT threshold. Using a higher threshold (e.g.,  $\geq$ 10) may lead to overly conservative labeling with potential false negatives, while lower thresholds increase coverage but may introduce noise. Thus, this temporal analysis provides critical insight into threshold selection and highlights the evolving nature of malware detection over time.

Figure 13 shows the t-SNE projections across varying VirusTotal (VT) detection (VirusTotal, 2025) thresholds. At lower thresholds (e.g., VT  $\geq$  4 to 6), there is significant overlap between the malware and benign clusters, indicating that many samples labeled as malware may exhibit similar characteristics to benign samples. This suggests that lower thresholds capture a broader range of potentially ambiguous or borderline malicious behaviors. As the VT detection count increases (e.g., VT  $\geq$  10), the overlap diminishes, and malware samples become more distinct and spatially separated from benign samples in the embedded space. This indicates that higher-threshold malware samples possess more distinguishable feature representations, likely reflecting stronger and more consistent malicious behaviors detected by a greater number of antivirus engines. Furthermore, the density of malware samples decreases as the threshold rises, aligning with the observed reduction in malware counts from the dataset.

To assess the impact of label noise on malware detection, we conduct a set of experiment varying the VirusTotal (VT) labeling Threshold. We first create a set of LAMDA datasets, where in each dataset contain *all benign samples* and *only malware samples with a specific VT labeling count*. For example, for the first dataset, we keep all benign samples and those malware samples that were flagged exactly by 4 antivirus engines (VT=4). In the next dataset, we include malware samples with VT=5, and so on, up to VT=13. This resulted in creating ten separate LAMDA dataset variants, each reflecting a different level of confidence in the AV engines.

Table 14: F1 scores of the baseline malware detectors with varying VT thresholds.

Split	Model	VT=4	VT=5	VT=6	VT=7	VT=8	VT=9	VT=10	VT=11	VT=12	VT=13
IID	LightGBM	0.8515	0.8566	0.8148	0.8397	0.8172	0.8511	0.8574	0.8364	0.6889	0.8385
	MLP	0.8364	0.8173	0.8166	0.8319	0.8393	0.8232	0.8102	0.8628	0.7536	0.8623
	SVM	0.7898	0.7954	0.7605	0.7543	0.8153	0.7492	0.7598	0.7353	0.6559	0.7462
	XGBoost	0.8456	0.8028	0.8298	0.8255	0.8119	0.8370	0.8273	0.7844	0.7324	0.7538
	LightGBM	0.2423	0.2008	0.2308	0.1994	0.2470	0.1962	0.1956	0.2178	0.1639	0.2133
NEAR	MLP	0.1971	0.1894	0.2191	0.2201	0.2223	0.2275	0.2249	0.1771	0.1867	0.2098
NEAK	SVM	0.1330	0.1512	0.1569	0.1763	0.1971	0.2712	0.1386	0.1381	0.1261	0.1865
	XGBoost	0.2221	0.2064	0.2703	0.2788	0.2761	0.1937	0.2788	0.2361	0.2563	0.2241
FAR	LightGBM	0.1284	0.1477	0.1559	0.1039	0.0852	0.1562	0.1551	0.0935	0.0511	0.2306
	MLP	0.1284	0.2716	0.0918	0.0721	0.1362	0.1873	0.1954	0.0841	0.0912	0.0947
	SVM	0.2393	0.2272	0.1503	0.0655	0.0938	0.1104	0.1322	0.0540	0.0508	0.0706
	XGBoost	0.1560	0.3513	0.2926	0.2613	0.1971	0.1787	0.2133	0.1067	0.1385	0.1452

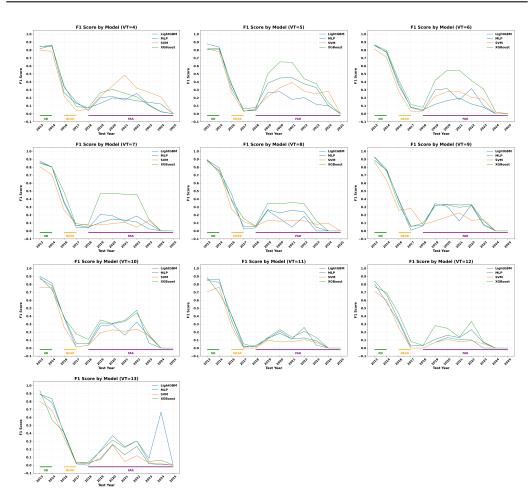


Figure 14: Combined F1 score plots for VT thresholds 4 to 13.

Next, We evaluate standard malware detectors (LightGBM, MLP, SVM, and XGBoost) on these datasets using the *AnoShift*-style splits, which simulate temporal concept drift by training on IID split and tested on NEAR and FAR splits. Each malware detector's performance is evaluated using F1-score metric.

Table 14 presents performance details of the baseline malware detectors using F1-scores metric with varying VT count. We made the following observations. LightGBM with VT=4, we observe an F1-score of 0.8515 under the IID split, however it drops significantly to 0.2423 on NEAR and to 0.1284 on FAR splits. This decline highlights the degradation of malware detector over time.

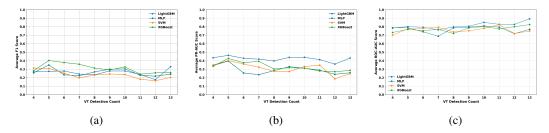


Figure 15: Comparison of model performance across VirusTotal detection thresholds for (a) F1-score, (b) PR-AUC, and (c) ROC-AUC.

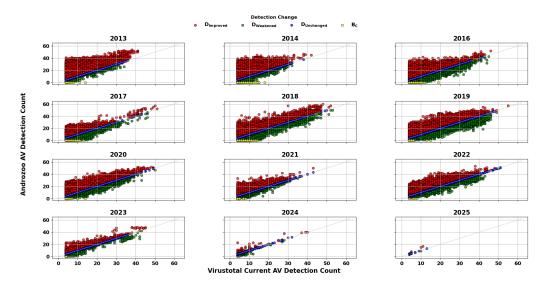


Figure 16: Malware AV detection drift over the years Virustotal vs AndroZoo Metadata.  $B_C$ : Currently Labeled as Benign,  $D_{Improved}$ : Improved Detection,  $D_{Weakened}$ : Weakened Detection,  $D_{Unchanged}$ : Unchanged Detection.

Figure 15 illustrates the average F1, PR-AUC, and ROC-AUC scores of the baseline malware detectors with varying VT labeling threshold. While Figure 15a is the summarization of Figure 14, subplots (a), (b) and (c) illustrate how performance metrics vary when the VT threshold is set to different values. Across all these three evaluation metrics, we observe only minor differences in baseline malware detectors performance. This suggests that, varying the VT labeling threshold has minimal impact on baseline accuracy. This result suggests that the primary causes of performance degradation in our main experiments may not be the labeling noise from VT, but rather factors such as temporal concept drift and class imbalance. The observed performance degradation is more likely attributed to the distributional shifts over time, reinforcing the relevance of concept drift in real-world malware detection scenarios.

## G ADDITIONAL VISUALIZATION ON LABEL DRIFT ACROSS YEARS

Figure 16 illustrates the analysis through a visual quantification of label drift across years, using Androozoo metadata and the updated VirusTotal (VT) (VirusTotal, 2025) report.

#### H SCALABILITY OF LAMDA

To support long-term use and extensibility, we have designed LAMDA with scalability in mind. In the context of LAMDA, scalability refers to the extensibility of the dataset—specifically, its ability to be easily expanded with new samples. We have published three variants of LAMDA on

the HuggingFace repository, each supporting a different VarianceThreshold configuration. The dataset creation process begins by splitting the static feature files (i.e., .data files extracted from each APK) into stratified train and test splits. From the training split, we collect the global set of all unique tokens (i.e., features), encode both train and test samples into binary vectors in this raw feature space, and apply VarianceThreshold to select high-variance features from the training data. The same selected features are then applied to the test data using the saved threshold object.

We publish the following artifacts to facilitate scalability: raw feature matrices (before thresholding), reduced feature matrices (after thresholding), and the serialized VarianceThreshold object (in joblib format). Using these resources and the accompanying codebase, researchers can seamlessly extend LAMDA by collecting newer APKs, extracting static features, encoding them, and applying the same thresholding object to map them into LAMDA's feature space. While it is not feasible to add new samples to the training set—because doing so would alter the global vocabulary and invalidate the original thresholding, researchers can add test-time samples for evaluation. This supports drift detection on newer and future malware variants without requiring retraining. Thus, LAMDA enables reproducible research and practical testing of detection models against evolving threats.

#### I CONTINUAL LEARNING ON LAMDA

In real-world settings, a large number of new benign and malicious Android applications are introduced each year. As a result, both benign (e.g., due to changes in user demands, Android APIs, security practices) and malicious (e.g., the emergence of novel malware variants) behaviors evolve over time, leading to concept drift. This makes it challenging for static machine learning models to maintain reliable performance over time without retraining regularly. However, complete retraining of past data becomes impractical due to the massive volume of Android applications released daily and the high computational cost associated with retraining. On the other hand, training solely on recent data often leads to catastrophic forgetting Lange et al. (2019); Rahman et al. (2022), where previously acquired knowledge is overwritten or lost. In such a situation, continual learning (CL) offers a compelling solution by enabling the detection models to adapt incrementally to new benign and malware applications without the need to retrain with all past data Park et al. (2025); Rahman et al. (2022); Ghiani et al. (2025). However, some CL techniques may require access to a small

# I.1 LAMDA FOR CONTINUAL LEARNING

subset of past data.

LAMDA can be a natural choice for benchmarking CL due to several key properties of its design and structure:

 I. **Temporal granularity:** It spans over a decade (2013–2025, excluding 2015) with available both monthly and yearly splits, allowing custom CL as per need.

II. Concept drift: As shown in Section 4, LAMDA exhibits significant distributional changes over time, both in feature and label space.

#### III. Flexible task construction:

 Domain-IL: Using yearly data splits while maintaining a consistent malware or benign labeling.

  Class-IL: Leveraging AVClass2-labeled malware families to incrementally expand the label space.

IV. Real world relevance: LAMDA is derived from real-world Android APKs and VirusTotal reports, introducing authentic drift and noise.

We evaluate CL on the LAMDA benchmark using two established baselines, *Naive* (i.e., None) and *Joint*, inspired by the prior work Rahman et al. (2022); Ghiani et al. (2025). Additionally, we include *Replay* (i.e., Experience Replay) Rolnick et al. (2019), a state-of-the-art memory replay based CL method, configured with a buffer size of 200 samples per experience. The Naive baseline trains the model sequentially on each experience or task without any access to past data, while the

Joint baseline retrains the model from scratch using the cumulative data observed up to the current experience or task.

These baselines are tested under two settings: Domain Incremental Learning (*Domain-IL*), which involves binary malware *vs* benign classification across yearly tasks, and Class Incremental Learning (*Class-IL*), where each task introduces new malware families to classify van de Ven et al. (2022); Rahman et al. (2022). Due to the lack of available prior work that can assign a single behavioral label, we didn't consider Task Incremental Learning (*Task-IL*) in our experimental setups.

We define each experience or task in the Domain-IL experiment as all samples (both benign and malicious) collected within a specific calendar year (e.g., 2013, 2014, ..., 2025). However, for the Class-IL experiments, each experience or task consists of only the malware samples collected during the corresponding year.

#### I.2 CONTINUAL LEARNING EXPERIMENTAL SETUP

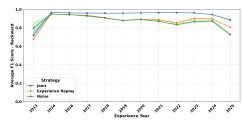
**Domain-IL.** In this setting, each experience or task corresponds to samples collected during a specific year (i.e., 2013, 2014, ... and so on). The model is designed to continuously learn to distinguish between malware and benign samples as the data distribution evolves over time. We use the Baseline variant of our published dataset, treating each year as a separate task in the learning sequence. The objective is for the model to adapt and maintain accurate binary classification performance despite the temporal distribution shifts.

Class-IL. In this setting, we utilize a different dataset derived from the Baseline variant of our published dataset. We selected only those malware families that contained more than 10 samples in the test set, resulting in a total of 154 families for our experiment. Consequently, we excluded the year 2025 from our experiments, as no family in that split met the minimum sample threshold. Additionally, we omit standard class incremental learning Park et al. (2025); Rahman et al. (2022), where entirely new classes are introduced in each experience. This approach does not reflect how malware appears in real-world scenarios, malicious samples often come from a mix of previously seen and new families. This claim is supported by the analysis presented in Table 5. The model is expected to learn incrementally to classify samples across all malware families encountered.

**Model Architecture.** We use a shared base architecture, a multi-layer perceptron (MLP) for both Class-IL and Domain-IL settings, consisting of four hidden layers — 512, 384, 256, 128, with ReLU activation. However, Task-specific heads are added to support each learning scenario. For Class-IL, we add a single linear layer outputting logits for all classes and train with categorical cross-entropy loss. For Domain-IL, we use a two-layer MLP head (100 units each, with dropout p=0.2) and a final sigmoid output, trained with binary cross-entropy. All networks are optimized with SGD (learning-rate 0.01, momentum 0.9, weight-decay 0.000001).

**Evaluation Metrics.** We evaluate classification performance using F1 score, which is the harmonic mean of precision and recall. It provides a balanced measure of a model's predictions, particularly important in *imbalanced* datasets. Following the prior work Ghiani et al. (2025), and we compute the F1 score after training on the k-th experience using two complementary evaluation modes:

- Backward Transfer Performance: We measures the model's ability to retain knowledge from previous tasks. After training on experience k, we compute the F1 scores on all previously seen experiences ( $\leq k$ ). This helps quantify the extent of *catastrophic forgetting* (CF).
- Forward Transfer Performance: We measures the model's ability to generalize to future, unseen tasks. After training on experience n, we compute the F1 score on all future experiences (>k). This indicates how well the model's current knowledge transfers to upcoming distributions.



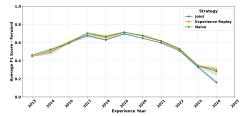
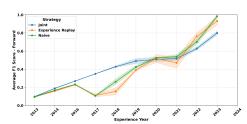


Figure 17: F1 Score in Domain-IL (Forward)

Figure 18: F1 Score in Domain-IL (Backward)



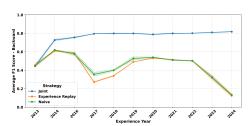


Figure 19: F1 Score in Class-IL (Forward)

Figure 20: F1 Score in Class-IL (Backward)

#### I.3 CONTINUAL LEARNING EXPERIMENTAL RESULTS

Figures 17, 18, 19, and 20 demonstrate the effectiveness of the CL methods in evaluating in realistic scenarios using LAMDA benchmark. In the Class-IL setting, we observe strong signs of catastrophic forgetting, especially in the Naive and Replay (Experience Replay) strategies. Backward F1 scores drop sharply after certain years, showing that learning new classes without retaining the previous knowledge leads to forgetting. Joint retains high performance as expected due to its exposure to all the previous data. In the Domain-IL setting, we observe that forgetting is relatively limited due to the fixed set of classes (malware or benign). Although the data distribution evolves over time, which leads to all strategies experiencing a gradual decline in forward F1 scores as they fail to adapt to new distributions. Additionally, we report the average F1 scores of LAMDA across all tasks under the Class-IL and Domain-IL scenarios in Tables 16 and 15. In the Class-IL setting (Table 16), the Joint strategy consistently achieves the highest performance, as expected, due to its access to the full dataset during training. However, this advantage also implies the need for significantly higher computational resources which makes it less practical for real-world settings. The Naive and Replay strategies perform considerably worse, which was also expected as the continual introduction of new classes. In contrast, the Domain-IL results (Table 15) show generally higher and more stable F1 scores across all strategies. Since the label space remains fixed over time, both Replay and even the Naive strategy perform reasonably well. This observation suggests that the primary challenge in Domain-IL is not always forgetting, but rather adapting to distributional shifts in the data.

These results highlight LAMDA's ability to capture both key challenges in continual learning: class expansion and distributional shift. As such, LAMDA serves as a realistic and challenging benchmark that supports future research in continual learning.

#### J COMPUTATIONAL RESOURCES FOR LAMDA GENERATION

All dataset processing and experiments for LAMDA were conducted on a high-performance compute server with the following configuration:

- CPU: Dual-socket Intel Xeon Gold 6430 with a total of 128 logical cores (64 physical cores, 2 threads per core).
- Memory: 1 TB RAM, with approximately 810 GB available during runtime.
- **GPU**: 4× NVIDIA H100 NVL GPUs with 95.8 GB memory per GPU. Experiments were conducted under CUDA 12.8 and driver version 570.124.06.

This infrastructure enabled us to efficiently process over 1 million APKs, large-scale temporal benchmarking over 12 years of Android malware data.

Table 15: Average F1 scores of Domain-IL across all experiences or tasks for LAMDA.

1464
1465
1466
1467
1468

Year Average F1 Score Average F1 Score Strategy Year Strategy Naive Naive  $48.86 \pm 1.15$  $80.26 \pm 0.27$  $83.08 \pm 0.13$  $79.93 \pm 0.27$ Joint  $46.89 \pm 0.33$  $46.56 \pm 0.00$ Joint Replay Replay

	Replay	40.50 ± 0.00	- 1	recpia	17.73 ± 0.21	
2014	Naive Joint Replay	$59.36 \pm 0.99$ $57.38 \pm 1.67$ $57.13 \pm 1.71$	2021	Naive Joint Replay	$78.61 \pm 0.33$ $84.12 \pm 0.18$ $79.46 \pm 0.55$	
2016	Naive Joint Replay	$68.07 \pm 0.82$ $69.21 \pm 0.78$ $68.75 \pm 1.20$	2022	Naive Joint Replay	$75.86 \pm 0.48$ $85.02 \pm 0.03$ $77.02 \pm 0.52$	
2017	Naive Joint Replay	$77.79 \pm 0.31$ $77.05 \pm 1.13$ $77.93 \pm 1.27$	2023	Naive Joint Replay	$78.10 \pm 1.28$ $85.59 \pm 0.23$ $80.74 \pm 0.73$	
2018	Naive Joint Replay	$76.64 \pm 1.06$ $76.61 \pm 0.15$ $75.94 \pm 1.46$	2024	Naive Joint Replay	82.54 ± 0.81 87.86 ± 0.11 84.77 ± 0.11	
2019	Naive Joint Replay	$79.74 \pm 0.16$ $82.76 \pm 0.41$ $79.72 \pm 0.09$	2025	Naive Joint Replay	72.71 ± 0.83 88.52 ± 0.46 80.57 ± 0.15	

Table 16: Average F1 scores of Class-IL across all experiences or tasks for LAMDA.

Year	Strategy	Average F1 Score	Year	Strategy	Average F1 Score
	Naive	12.60 ± 0.65	2020	Naive	$53.50 \pm 0.18$
2013	Joint	$12.89 \pm 0.37$		Joint	$68.63 \pm 0.67$
	Replay	$12.75 \pm 0.48$		Replay	$52.67 \pm 1.34$
2014	Naive	25.01 ± 0.28	2021	Naive	51.83 ± 0.87
	Joint	$28.64 \pm 0.49$		Joint	$72.23 \pm 0.46$
	Replay	$24.07 \pm 0.74$		Replay	$50.21 \pm 1.51$
2016	Naive	32.64 ± 0.72	2022	Naive	53.91 ± 0.42
	Joint	$40.14 \pm 0.22$		Joint	$76.98 \pm 0.36$
	Replay	$32.60 \pm 0.36$		Replay	$54.74 \pm 0.43$
2017	Naive	19.68 ± 1.12	2023	Naive	38.39 ± 1.84
	Joint	$51.14 \pm 0.09$		Joint	$80.88 \pm 0.25$
	Replay	$16.92 \pm 1.21$		Replay	$37.65 \pm 2.46$
2018	Naive	32.47 ± 2.06	2024	Naive	13.27 ± 1.33
	Joint	$59.66 \pm 0.31$		Joint	$81.79 \pm 0.29$
	Replay	$23.86 \pm 2.22$		Replay	$12.49 \pm 1.76$
2019	Naive	47.98 ± 1.71			
	Joint	65.86 ± 1.20			
	Replay	$44.43 \pm 0.26$			