# TRAMEL: AN EXEMPLAR REPLAY-BASED CONTIN-UAL LEARNING FRAMEWORK FOR MALWARE TRAFFIC ANALYSIS

**Anonymous authors**Paper under double-blind review

# **ABSTRACT**

Most prior work on continual malware detection has centered on static code analysis. In contrast, this paper explores continual learning (CL) for malware traffic analysis, which leverages encrypted flow features to capture behavioral signals resilient to obfuscation and encryption. Unlike intrusion detection systems that focus on coarse anomaly detection, malware traffic analysis requires fine-grained family-level classification under evolving and imbalanced distributions, making it a distinct and challenging setting for CL.

We introduce TraMEL (Traffic-based Malware Exemplar Learning), a replay-based CL framework designed for malware traffic. TraMEL integrates (i) adaptive exemplar selection to balance long-tailed family distributions and (ii) an exemplar refinement phase to counter task recency bias while operating under strict memory budgets. We evaluate both standard class-incremental and temporally shifted scenarios, showing that TraMEL consistently outperforms strong CL baselines such as iCaRL, ER, and TAMiL by 10–30 percentage points on CI-CAndMal2017 and IoT23. Remarkably, it approaches the performance of joint training, a theoretical upper bound with full access to past data. These results demonstrate that CL on malware traffic is both feasible and practical, providing a memory-efficient approach for real-world detection. Code is available at https://anonymous.4open.science/r/ICLR2026-code-D575/.

# 1 Introduction

Modern malware increasingly evades traditional defenses by using encryption (e.g., TLS 1.3) to render packet inspection ineffective and code obfuscation to bypass static analysis (Moser et al., 2007; Deng & Mirkovic, 2022; Anderson & McGrew, 2016). This shift motivates malware traffic analysis (MTA): detecting malicious activity directly from encrypted network traffic flows rather than executable code. Unlike conventional intrusion detection systems (IDS), which rely on coarse logs or binary anomaly flags under closed-world assumptions (Sommer & Paxson, 2010; Paya et al., 2024), MTA requires fine-grained family-level classification in an open world where malware families continually evolve and reappear (Mariconti et al., 2017). This setting introduces distinct learning challenges, especially the need to adapt without retraining from scratch (Rahman et al., 2022). The challenge is particularly acute in mobile and embedded ecosystems, where encrypted traffic dominates and malware is fast-changing. To capture this, we study two representative domains: Android malware, using the CICAndMal2017 dataset with 42 families (Lashkari et al., 2018), and IoT malware, using the IoT-23 dataset featuring botnets such as Mirai (Garcia et al., 2020).

Machine learning (ML) models have achieved strong results on static malware traffic analysis (MTA) benchmarks (Mirsky et al., 2018; Anderson & McGrew, 2016). We argue, however, that this success reflects a flawed *closed-world assumption*, which does not hold in real, evolving networks (Sommer & Paxson, 2010). In practice, network traffic is highly non-stationary: *distributional drift* arises as malware changes its command-and-control infrastructure, protocol variation, and benign applications modify their communication behavior (Küchler et al., 2021). More importantly, this drift is often the result of *adversary-induced non-stationarity*, where attackers deliberately generate variants of known malware and introduce new families to bypass detection (see Section 2). Such dynamics steadily erode classifier performance. The standard countermeasure—fine-tuning on new

055

057

060

061

062

063 064

065

066 067

068

069

070

071

072

073 074

075

076

077

079

081

082

083

084

085

087

880

089

090

091

092

093

094

095 096

097

098

099

100

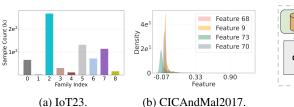
101

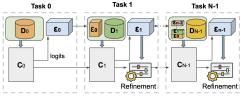
102

103 104

105 106

107





skewed feature distribution in CIC2017 (1b).

Figure 1: Class imbalance in IoT23 (1a) and Figure 2: Overview of the proposed TraMEL pipeline.

data—proves unreliable, as it triggers catastrophic forgetting (CF), where the model loses its ability to detect previously learned malware families while adjusting to new threats.

To the best of our knowledge, no malware traffic dataset captures the long-term recurrence of families, making it a hard problem to directly evaluate forgetting under realistic conditions. To address this gap, we design two benchmarks. The first is a standard class-incremental (Class-IL) split with disjoint families, representing a strict worst-case. The second is a temporal split based Class-IL, grouping families by time of first appearance to mirror natural traffic shifts. Although both exclude recurring families, the temporal split is deliberately conservative, harder than real deployments, thus providing a principled benchmark for continual learning (CL) in MTA.

As such, in this work, we formalize MTA as a Class-IL problem with three objectives: (i) preserve accuracy on past families, (ii) adapt to emerging ones, and (iii) operate under strict memory and compute budgets. Memory replay is particularly well-suited here, as it retains knowledge through compact exemplar buffers (Rahman et al., 2025). While methods such as ER (Rolnick et al., 2019), iCaRL (Rebuffi et al., 2017a), and TAMiL (Bhat et al., 2023b) have proven effective in vision, prior CL work in IDS settings (Channappayya et al., 2023; Amalapuram et al., 2024) has focused on coarse binary anomaly detection under closed-world assumptions. These approaches do not address the fine-grained, long-tailed, and evolving family distributions characteristic of malware traffic. While prior studies propose CL for malware analysis to cope with temporal drift (Sun et al., 2025; Park et al., 2025; Rahman et al., 2025), but these efforts focus on code-level analysis. In contrast, our work addresses encrypted MTA, tackling catastrophic forgetting across both classincremental and temporal shift scenarios through exemplar replay and refinement.

Our approach. We introduce TraMEL (Traffic-based Malware Exemplar Learning), an exemplarreplay CL framework specifically tailored for MTA. TraMEL is designed to tackle three core challenges in this domain. 

Real-world malware traffic exhibits long-tailed family distributions and sparse feature vectors (Figure 1). To cope with this imbalance, TraMEL employs exemplar selection strategies that balance class representation while capturing intra-class diversity. 2 Models trained incrementally are prone to task recency bias, where knowledge of earlier families is degraded as new ones are introduced. To mitigate this effect, TraMEL introduces a refinement phase that fine-tunes exclusively on buffered exemplars, reinforcing prior knowledge without revisiting the full history. 3 Finally, malware detectors must operate under strict memory budgets. TraMEL therefore emphasizes compact but representative exemplar selection, ensuring long-term accuracy even under severe buffer constraints.

To this end, TraMEL combines a heuristic exemplar selection strategy—balancing class coverage with diversity-aware clustering—with an exemplar refinement phase that replays buffered samples to mitigate forgetting while maintaining adaptability.

Results. On CICAndMal2017 and IoT23, TraMEL consistently outperforms strong CL baselines such as iCaRL, ER, and TAMiL. Even with a buffer of only 3,000 samples (0.5% of data), it achieves 10-30 percentage points higher accuracy and approaches joint-level results. Clustering-based selection is especially effective under tight memory, while simpler strategies suffice with larger buffers.

# THREAT MODEL

Retrograde Malware Attack (RMA) targets ML-based malware detectors that are incrementally updated with only *new* traffic or file samples. In practice, security pipelines often retrain classifiers on fresh threat intelligence feeds (e.g., new flows, domains, binaries) without retaining historical corpora due to storage and scalability limits. This induces *catastrophic forgetting* of earlier malware signatures and behavioral traces, enabling adversaries to weaponize legacy or lightly modified variants that evade detection. From the perspective of network traffic analysis,

RMA unfolds in three phases (Figure 3):

- Initial Training (1): The detector is trained on malicious and benign traffic (e.g., packet sequences, flow metadata, TLS fingerprints).
- Updates and Forgetting (2): The model is periodically retrained on recent captures (e.g., from honeypots or sandboxes). Because older families and domains are excluded, recall on previously seen traffic patterns declines while benign software may be misclassified, raising false positives.

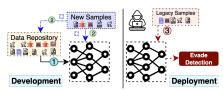


Figure 3: Retrograde Malware Attack (RMA).

• RMA in Deployment (3): Adversaries exploit this forgetting by reintroducing legacy families or slightly altered variants.

# 3 RELATED WORK

**Replay in CL.** Addressing CF is the core challenge in CL, and one widely used solution is <u>replay</u>. Replay methods improve learning by mixing current data with representative information from earlier tasks. They are typically grouped into two categories – exact replay (storing real samples) and generative replay (generating synthetic data). Exact replay stores a fixed number of past samples, controlled by a memory budget  $\mathcal{M}$ . Methods like ER (Rolnick et al., 2019), A-GEM, and iCaRL aim to maintain performance while using as few replay samples as possible (Rolnick et al., 2019; Chaudhry et al., 2019; Rebuffi et al., 2017a). TAMiL (Bhat et al., 2023a) builds on ER by using attention to retain prior data distributions at the representation level, improving knowledge retention beyond simple replay. Generative or pseudo-replay strategies are designed to replicate the original data (Li & Hoiem, 2017; Shin et al., 2017; van de Ven et al., 2020). These techniques either generate a representative of the original data using a separate generative model or generate pseudo-data by using an earlier model's predictions as soft labels for training subsequent models.

CL in Malware and Related Domains. Study of CL in malware domains is relatively limited. Rahman et al., Rahman et al., 2022) showed that replay-based methods are more effective due to the structured and diverse nature of tabular malware features. MalCL(Park et al., 2025) extends this with a GAN-based generative replay and feature-guided sampling, while MADAR (Rahman et al., 2025) introduces distribution-aware replay to select representative and discriminative samples. Beyond malware classification, other efforts address adjacent problems. Chen et al., 2023) study concept drift in Android malware using contrastive and active learning, but do not tackle CF.

SPIDER (Amalapuram et al., 2024) extends CL to intrusion detection using a semi-supervised approach that matches supervised baselines while storing only unlabeled traffic. Still, it operates under a closed-world binary setting and requires up to 20% labeled data—limiting its practicality for malware. Its companion, Augmented-Memory Replay (Channappayya et al., 2023), uses only intrusion benchmarks, with limited relevance to real malware traffic and no support for privacy-preserving replay. Other work frames CL in the context of network traffic but still falls short. SPCIL (Xu et al., 2024) introduces a lightweight dual-branch model for malware detection but handles only small class increments, with growing memory and stability concerns. Zhang et al. (Zhang et al., 2025) propose an expandable CL system with per-task frozen extractors and neural architecture search. While effective on IoT and VPN datasets, these settings lack family-level malware structure and do not scale to long-horizon, evolving malware detection.

Current CL systems for network security either frame intrusion detection as a binary anomaly task or evaluate on IoT/VPN traffic, which lacks the family-level diversity characteristic of real malware. Notably, existing work does not address CL for discovering and adapting to new malware families. These limitations motivate our focus on TRAMEL, which directly targets CF in the context of evolving malware families and realistic traffic streams.

# 4 OVERVIEW OF TRAMEL

We present TraMEL, a continual learning framework for malware traffic classification in a class-incremental (Class-IL) setting. We assume that new malware families (i.e., classes) arrive incrementally over tasks  $t_0, \ldots, t_{n-1}$ , each associated with a training set  $D_0, \ldots, D_{n-1}$ . At task  $t_i$ , the objective is to train a classifier  $C_i$  on the current dataset  $D_i$  while retaining knowledge from previous datasets  $D_0, \ldots, D_{i-1}$ .

In our experiments, we consider both synthetic Class-IL splits and more realistic temporal shifts. In the Class-IL setup, tasks are defined by evenly partitioning malware families across n tasks, which stresses the ability to recognize new families while preserving old ones. In the temporal setup, tasks are organized by the year in which malware families first appear, mimicking how new variants emerge in practice. This allows us to evaluate TraMEL under conditions where distributions evolve naturally over time, reflecting adversary-driven drift.

TraMEL addresses these scenarios through a three-phase process. First, the model is trained jointly on the current task data  $D_i$  and the replay buffer  $E_{< i}$  containing exemplars from earlier tasks, reducing early forgetting. Second, a set of informative exemplars is selected from  $D_i$  under a fixed memory budget. The selection strategy explicitly promotes class balance and intra-class diversity, ensuring that even minority families are preserved in the buffer. Finally, to mitigate task recency bias, the model is refined exclusively on the buffer  $E_i$ , consolidating older knowledge without requiring full historical data.

By combining joint training, imbalance-aware exemplar selection, and targeted refinement, TraMEL achieves a balance between plasticity and stability across both synthetic and temporally defined tasks. This enables robust long-term malware detection in evolving threats. The following subsections detail the exemplar selection strategy, buffer management, and refinement procedure.

#### 4.1 Exemplar Selection Strategies

Let  $\mathcal{D}^c = \{(x_j, y_j)\}_{j=1}^{N_c}$  denote the training samples of class c, and let f(x) be the feature representation of input x extracted by the backbone network. The goal is to select m exemplars  $\mathcal{E}^c \subset \mathcal{D}^c$  for each class to be stored in the replay buffer. We investigate three strategies.

The first is random sampling, which simply draws m samples uniformly from  $\mathcal{D}^c$ . This baseline provides unbiased coverage of the class distribution but does not exploit structure in the feature space.

The second is *class-mean selection*, following iCaRL (Rebuffi et al., 2017b). We compute the class prototype

$$\mu_c = \frac{1}{N_c} \sum_{(x_j, y_j) \in \mathcal{D}^c} f(x_j),$$

and select the m samples with the highest similarity to  $\mu_c$ . This aligns exemplars with the class centroid, ensuring representativeness, though it may suffer from limited diversity.

The third is a *clustering-based strategy* to enhance diversity. Feature vectors  $\{f(x_j)\}$  are partitioned into k clusters using K-means with Euclidean distance, yielding centroids  $\{\mu_{c,1}, \ldots, \mu_{c,k}\}$ . Each cluster i receives a quota  $m_i$  proportional to its size:

$$m_i = \left\lfloor m \cdot \frac{|\mathcal{D}_i^c|}{N_c} \right\rfloor.$$

From each cluster, we select the  $m_i$  samples closest to its centroid:

$$\mathcal{E}_{i}^{c} = \arg \min_{\substack{\mathcal{S} \subset \mathcal{D}_{i}^{c} \\ |\mathcal{S}| = m_{i}}} \sum_{x \in \mathcal{S}} \|f(x) - \mu_{c,i}\|_{2}^{2}.$$

The final exemplar set is  $\mathcal{E}^c = \bigcup_{i=1}^k \mathcal{E}^c_i$ . By enforcing coverage of multiple clusters, this method captures diverse semantic regions, mitigating over-representation of dense areas and improving generalization under continual learning. We empirically find that using larger numbers of clusters (e.g.,  $k \geq 100$ ) further improves performance on CICAndMal2017 and IoT23, as the buffer more faithfully reflects the underlying data manifold. Detailed results are provided in the Appendix  $\ref{eq:continuous}$ ?

#### 4.2 Replay Buffer

Storing all past data for retraining is infeasible; instead, TraMEL maintains a fixed-size replay buffer of capacity K to hold exemplars from earlier tasks. In the Class-IL setting, the number of classes grows over time while K remains constant, so the quota per class decreases as tasks accumulate. If  $M_i$  denotes the number of classes introduced at task i, then after task i each class receives  $\frac{K}{\sum_{j=1}^i M_j}$  exemplars. This progressive reduction makes exemplar quality increasingly critical.

Compared to vision benchmarks, where  $K \leq 1,000$  (roughly 3% of training data) (Rebuffi et al., 2017b), malware traffic datasets require much larger buffers due to their scale. For example, maintaining the same ratio on CICAndMal2017 implies  $K \approx 33,000$ . Such scale exacerbates memory constraints and highlights the need for selection strategies that emphasize both representativeness and diversity.

To capture these practical considerations, we evaluate TraMEL under multiple buffer capacities proportionally scaled to dataset size (from 200 to 60,000), enabling a systematic analysis of how memory budgets influence exemplar effectiveness.

# 4.3 EXEMPLAR REFINEMENT

In the *i*-th task, training on the current dataset  $D_i$  together with the exemplar buffer  $E_{< i}$  creates a severe imbalance, since  $|D_i| \gg |E_{< i}|$ . This imbalance amplifies CF and leads to task recency bias, where the model favors recently observed classes (Lyu et al., 2023).

To counter this effect, TraMEL introduces a refinement phase after each task. In this phase, the model is fine-tuned exclusively on the exemplar buffer  $E = E_{< i} \cup E_i$ , which acts as a compact proxy for past distributions. Since exemplars are carefully selected for both representativeness and diversity, replaying them provides an efficient rehearsal step.

The refinement objective integrates supervised and distillation losses to balance plasticity and stability. Let  $f^{(i)}(x)$  be the logits of the current model after refinement on task i,  $f^{(i-1)}(x)$  the logits from the previous refined model, and  $f^{(i)'}(x)$  the logits from the model immediately after task i training. For an exemplar x with label y, we define:

$$\mathcal{L}_{\text{refine}} = \mathcal{L}_{\text{CE}} + \alpha \cdot \mathcal{L}_{\text{past}} + \beta \cdot \mathcal{L}_{\text{current}},$$

where

$$\mathcal{L}_{\text{CE}} = \frac{1}{|E|} \sum_{(x,y) \in E} \text{CE}(f^{(i)}(x), y), \quad \mathcal{L}_{\text{past}} = \frac{1}{|E_{

$$\mathcal{L}_{\text{current}} = \frac{1}{|E_i|} \sum_{x \in E_i} \|f^{(i)}(x) - f^{(i)'}(x)\|_2^2.$$$$

Here,  $\mathcal{L}_{CE}$  enforces correct classification across all exemplars,  $\mathcal{L}_{past}$  preserves behavior on earlier tasks by aligning with the previous refined model, and  $\mathcal{L}_{current}$  stabilizes adaptation to the new task by constraining deviation from the post-training model. Together, these terms mitigate recency bias while preventing overcorrection, yielding a refined balance between adaptation to emerging malware families and retention of prior knowledge. Hyperparameters  $\alpha$ ,  $\beta$ , and the number of refinement epochs are scaled with buffer size and task composition; detailed sensitivity analyses are reported in Section 5.4.

## 4.4 Classifier Architecture

Malware traffic data is inherently tabular, with each flow represented by dozens of statistical and protocol-level features rather than raw sequences or images. To identify a suitable backbone for continual learning, we evaluate three neural architectures: Multi-Layer Perceptrons (MLPs), one-dimensional Convolutional Neural Networks (CNNs), and Vision Transformers (ViTs).

The MLP baseline consists of nine fully connected layers with ELU/ReLU activations, batch normalization, and dropout. While computationally efficient, it provides limited representational power and yields the weakest performance. The CNN baseline uses a six-layer 1D convolutional stack with

max pooling and fully connected layers (28M parameters). Although it achieves up to 75% accuracy, it suffers from instability across runs and sensitivity to initialization, reflecting the difficulty of applying local convolutional filters to tabular features without strong positional structure.

In contrast, the Transformer-based model delivers both higher accuracy and greater stability. On CICAndMal2017, a ViT with six encoder blocks (hidden dimension 384, MLP size 1152, eight heads) achieves 75–80% accuracy with only 8.9M parameters. On IoT23, a lighter configuration (hidden size 16, MLP size 48, one encoder layer, two heads) achieves competitive accuracy despite the smaller input dimension. In both cases, the ViT consistently outperforms CNNs and MLPs in average accuracy and variance, while maintaining robustness across the entire Class-IL sequence.

These results provide an important insight: attention-based models are particularly well-suited for malware traffic analysis. Unlike CNNs, which rely on local receptive fields, Transformers capture global inter-feature dependencies without assuming positional priors, making them effective on tabular data where relationships among features (e.g., packet size, timing, DNS queries) are long-range and non-sequential. Moreover, the ViT achieves stronger accuracy—complexity tradeoffs, with fewer parameters yet higher stability than CNNs. This aligns with recent evidence that Transformers generalize well to structured tabular data (Huang et al., 2020).

## 5 EXPERIMENTAL DETAILS

## 5.1 Dataset

We evaluate TraMEL and other replay-based CL models on two publicly available malware traffic datasets: CICAndMal2017 (Lashkari et al., 2018) and IoT23 (Garcia et al., 2020). Both datasets are split into training, validation, and test sets using an 8:1:1 ratio.

**CICAndMal2017 (1,105,290 flows).** This dataset consists of Android malware traffic spanning 42 families. During preprocessing, IP and port fields are anonymized, and traffic direction is inferred using a manually defined list of local IP addresses. Timestamps are normalized to compute interpacket delays (IPD), which are further adjusted by traffic direction. The dataset is highly imbalanced, with the largest family containing over 75,000 flows and the smallest fewer than 4,000.

**IoT23** (712,231 flows). This dataset contains IoT network traffic of 11 malware families. To improve class balance, we exclude two minority families (Torri and Trojan), resulting in 9 classes. Preprocessing removes timestamps, unique identifiers, host addresses, and tunneling or service-related fields. Numeric packet and byte features are log-transformed to reduce skewness. Similar to CICAndMal2017, the class distribution is highly imbalanced Hideandseek (~267k flows), Linux.Hajime (131k), and Muhstik (114k) dominate, while families like Hakai contain as few as 4,000 flows.

#### 5.2 TASK CONFIGURATION AND TRAINING PROTOCOL

We evaluate two class distribution scenarios in a Class-IL setting. The first follows prior work showing that assigning more classes to the initial task can mitigate forgetting in subsequent tasks (Park et al., 2025; Rahman et al., 2022). For CICAndMal2017, we configure tasks as  $M_1=22$  and  $M_2=M_3=M_4=M_5=5$ . For IoT23, we set  $M_1=5$  and  $M_2=M_3=M_4=M_5=1$ .

The second scenario is motivated by the fact that malware families often reappear over time as new variants (Sun et al., 2025). To the best of our knowledge, no publicly available malware traffic dataset captures the same families re-emerging over time, which prevents a direct evaluation of how temporal evolution affects malware traffic analysis. Since our datasets were collected over relatively short periods, we approximate temporal dynamics by grouping malware families according to their time (year) of first appearance. For CICAndMal2017, this yields  $M_1=4,\,M_2=5,\,M_3=7,\,M_4=4,\,M_5=10,\,M_6=6,$  and  $M_7=6.$ 

It is worth noting that in our setting, all families across tasks are disjoint. This makes the split stricter than real deployments, where families may persist and reappear, enabling transfer. Nevertheless, the overall malware-traffic distribution still shifts across tasks, so the setup remains meaningful for assessing distributional non-stationarity. Our results should be viewed as a conservative lower bound; in practice, temporal reoccurrence would likely ease the problem. Each experiment is repeated five

times and we report mean accuracy; training uses 50 epochs on CICAndMal2017 and 40 on IoT23 with early stopping after the first task.

In the refinement phase, we fix a constant k to balance buffer size and refinement epochs, ensuring consistent replay across settings. For CICAndMal2017, k=240,000, and for IoT23, k=20,000. This value is determined empirically and scales with buffer size and dataset scale.

# 5.3 EVALUATION METRICS.

We report task-wise and mean accuracy as primary metrics. For each task i, accuracy is computed over all test samples from classes seen up to i, capturing both new learning and retention. Mean accuracy is the average of task-wise results, reflecting overall stability and forward transfer. To quantify CF, we use the forgetting score, defined as the per-class gap between maximum and current accuracy, which also serves to assess task recency bias.

#### 5.4 Hyperparameter Tuning

We tune three key hyperparameters on CICAndMal2017 (buffer size, refinement epochs, and distillation weights  $(\alpha, \beta)$  and evaluate their impact using mean accuracy and forgetting score.

**Buffer size and refinement epochs.** We vary buffer sizes between 3,000 and 33,000 and adjust refinement epochs (80 vs. 8) to keep the total number of exemplar updates per task fixed at 240K. As shown in Table 5, increasing refinement epochs effectively compensates for smaller buffers, improving retention of past knowledge.

**Distillation weights**  $(\alpha, \beta)$ . We tune  $\alpha$  to preserve past-task knowledge and  $\beta$  to emphasize current-task accuracy. While  $\alpha = \beta = 1$  already stabilizes learning, unbalanced settings reveal a trade-off: larger  $\alpha$  improves retention but reduces new-task accuracy, whereas larger  $\beta$  favors recent tasks at the cost of earlier ones. As shown in Table 5, mean accuracy remains similar across settings, but forgetting scores vary significantly, highlighting the importance of tuning  $(\alpha, \beta)$  for stability-plasticity balance.

#### 6 RESULTS

Comparison to Baselines. We evaluate TraMEL against replay-based CL methods including iCaRL (Rebuffi et al., 2017b), ER (Rolnick et al., 2019), and TAMiL (Bhat et al., 2023b) on CI-CAndMal2017 and IoT23, using the same buffer size of K=33,000 exemplars. For reference, we also report two standard baselines: *None*, which trains only on the current task without replay, and *Joint*, which serves as an accuracy upper bound with full access to past and current data. Table 1 shows that TraMEL consistently outperforms iCaRL, ER, and TAMiL in both mean accuracy and stability, achieving 10–30 percentage points higher accuracy on the final task. Results in Table 2 highlight the effect of refinement: although TraMEL trails TAMiL in the earliest tasks, it surpasses all baselines in later tasks under temporal drift, demonstrating stronger retention and reduced recency bias.

Compared to the Joint baseline, TraMEL narrows the gap substantially, by about 13 percentage points on CICAndMal2017 (including the temporal shift scenario) and nearly matching Joint performance on IoT23. In terms of efficiency, training with  $K=33{,}000$  exemplars remains practical. On an NVIDIA A100 80GB PCIe GPU, the Joint baseline required about 2.5 hours, whereas TraMEL-K completed in roughly one hour, achieving competitive accuracy with significantly lower computational cost.

**Exemplar Selection.** We adopt random sampling (TraMEL-R), centroid-based selection (C-mean), K-means clustering with 600 clusters (TraMEL-K) on IoT23. TraMEL-R proves effective when memory is sufficient, while TraMEL-K provides greater robustness under tighter memory budgets by enhancing exemplar diversity. We further analyze the effect of varying the number of clusters in the Appendix A.2, which shows that larger k values generally improve coverage of long-tailed distributions, peaking at around  $N_k = 600$ , and slightly degrades beyond that point.

378 379

Table 1: Performance of TraMEL on CICAndMal2017 (CIC17) and IoT23 datasets.

3	ξ	3	(
3	8	3	1
3	8	3	6
3	8	3	3
3	8	3	2
3	8	3	Į,
3	8	3	6
3	8	3	1
3	ξ	3	ξ

38	36
38	37
38	88
38	39
39	0
39	1
39	)2

393 394 395

396 397 398 399 400

404 405 406 407

408 409

410 411 412

413

414

415

416

417

418 419 420 421 422 423 424

425

426

427

428 429

430

431

Dataset | Model Task 1 Task 2 Task 3 Task 4 Task 5 Mean  $77.12 \pm 2.5$  $75.71 \pm 1.8$  $76.00 \pm 1.4$  $76.25 \pm 0.9$  $75.61 \pm 0.3$  $76.14 \pm 1.1$ Joint None  $77.12 \pm 2.5$  $14.92 \pm 2.6$  $14.84 \pm 3.1$  $13.37 \pm 3.6$  $10.55 \pm 2.7$  $26.16 \pm 1.4$ CIC17 TraMEL-R  $76.09 \pm 1.9$  $66.54 \pm 2.0$ **60.36**  $\pm$  1.4  $56.63 \pm 1.3$  $53.75 \pm 0.6$  $62.67 \pm 1.2$ ER  $55.23 \pm 24.3$  $55.75 \pm 19.4$  $59.28 \pm 4.7$  $54.18 \pm 2.9$  $39.60 \pm 18.7$  $52.81 \pm 9.5$  $33.25\pm6.5\phantom{0}$ iCaRL  $55.16 \pm 6.0$  $29.59 \pm 17.5$  $30.78 \pm 11.8$  $28.31 \pm 4.8$  $22.39 \pm 2.6$ **TAMiL**  $57.69 \pm 15.9$  $56.18 \pm 10.8$  $48.79 \pm 18.5$  $31.44 \pm 25.4$  $47.15 \pm 7.4$  $48.25 \pm 6.3$  $89.55 \pm 8.6$  $88.11 \pm 8.5$  $85.67 \pm 5.5$  $82.15 \pm 4.1$  $81.99 \pm 1.0$  $85.50 \pm 4.3$ Joint  $23.92\pm24.6\phantom{0}$  $89.55 \pm 8.6$  $6.69 \pm 7.1$  $12.53 \pm 14.4$  $29.57 \pm 7.1$ None  $15.14 \pm 8.6$ IoT23 TraMEL-K  $89.54 \pm 9.0$  $82.65 \pm 9.0$  $76.34 \pm 10.0$  $63.07 \pm 11.0$  $59.17 \pm 18.0$  $74.15 \pm 10.0$  $43.19 \pm 15.1$ iCaRL  $67.37 \pm 22.7$  $67.93 \pm 18.2$  $65.49 \pm 9.2\phantom{0}$  $59.70\pm7.1$  $54.51 \pm 15.2$ ER  $78.52 \pm 13.3$  $88.21 \pm 10.9$  $70.11 \pm 12.0$  $70.17 \pm 8.8$  $54.29 \pm 22.1$  $72.26 \pm 2.8$  $51.06\pm14.8$  $59.30 \pm 13.8$ **TAMiL**  $\mathbf{81.23} \pm \mathbf{18.1}$  $64.20 \pm 11.0$  $47.51 \pm 18.6$  $52.51 \pm 15.8$ 

Table 2: Performance of TraMEL on CICAndMal2017 in the temporal drift setting.

Model	Task 1	Task 2	Task 3	Task 4	Task 5	Task 6	Task 7	Mean
Joint None	$78.71 \pm 1.7$ $78.49 \pm 0.8$				$69.32 \pm 0.4$ $31.72 \pm 0.1$	l		
TraMEL-R	78.16 ± 2.0	I	I		<b>54.69</b> ± 1.1	I		<u> </u>
ER	$79.18 \pm 14.3$	$40.55 \pm 25.8$	$53.87 \pm 1.4$	$55.32 \pm 2.2$	$46.83 \pm 0.4$	$41.17 \pm 7.7$	$43.66 \pm 1.6$	$51.51 \pm 4.8$
iCaRL	$71.47 \pm 7.6$	$54.60 \pm 5.5$	$43.53 \pm 3.6$	$33.95 \pm 7.7$	$14.33 \pm 8.0$	$22.56 \pm 3.3$	$22.06 \pm 3.3$	$37.5 \pm 5.6$
TAMiL	$82.44 \pm 3.7$	<b>65.67</b> $\pm$ 1.6	$49.04 \pm 9.8$	$44.98 \pm 5.4$	$36.45 \pm 17.8$	$38.26 \pm 6.4$	$38.28 \pm 14.9$	$50.73 \pm 5.1$

Table 3: Performance of TraMEL on the IoT23 dataset with different exemplar selection strategies.

Method	Task1	Task2	Task3	Task4	Task5	Mean
Random	98.38	90.54	79.94	69.91	73.39	82.43
C-Mean	98.38	88.46	79.74	65.19	72.18	80.79
$KM(N_k = 600)$	98.38	90.78	80.40	69.28	75.08	82.78

**Replay Buffer Size.** To examine how buffer size influences CL performance, Table 4 reports results on CICAndMal2017 across seven capacities: K=200,500,1,000,3,000,6,000,33,000, and 60,000. All methods improve with larger buffers, but ER performs best at very small sizes (K<1,000). Once the buffer reaches 1,000 exemplars, TraMEL consistently achieves higher accuracy, exceeding baselines by more than 10 percentage points when the buffer is sufficiently large to represent each class.

We also observe differences between TraMEL-R (random sampling) and TraMEL-K (K-means selection). Under tight memory budgets (K=200–6,000), TraMEL-K performs better by enforcing greater exemplar diversity. As buffer size increases, random sampling becomes adequate to capture representative samples, and the performance gap between the two strategies narrows.

**Exemplar Refinement.** A key challenge in refinement is balancing adaptation to new tasks with retention of prior knowledge. This is especially critical in malware classification, where detecting newly emerging families must not come at the expense of forgetting earlier ones. Figure 4 illustrates this effect at the task level: before refinement, predictions are skewed toward Task 5, reflecting severe recency bias; after refinement, they are more evenly distributed, indicating improved stability. For detailed class-level confusion matrices, see Appendix A.1.

Figure 5 further shows how the refinement loss weights  $\alpha$  and  $\beta$  shape this trade-off. With  $\alpha=4,\beta=1$  (Figure 5a), earlier tasks improve by over 10 percentage points, though Task 5 drops by 15% – still within 5% of the Joint baseline. Conversely, with  $\alpha=1,\beta=4$  (Figure 5b), Task 5 is better preserved but forgetting of earlier tasks is more severe.

Table 4: Mean accuracy under varying buffer sizes on CICAndMal2017.

Method	200	500	1K	3K	6K	33K	60K
TraMEL-R	30.83	34.16	38.37	47.12	52.69	64.20	66.42
TraMEL-K	31.13	35.32	39.20	49.15	54.59	61.88	63.44
iCaRL	22.19	27.30	27.88	32.19	32.66	30.76	32.54
TAMiL	27.43	32.90	30.21	36.17	44.86	44.24	50.03
ER	35.61	37.15	39.28	36.53	44.55	51.42	57.11

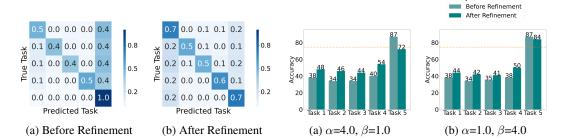


Figure 4: Task-level normalized confusion matrix on CIC17 (Task 5). *Before* refinement, predictions are biased toward the current task; *after*, they are more evenly distributed, indicating reduced recency bias and forgetting.

Figure 5: Per-task accuracy on CIC17 after refinement. Larger  $\alpha$  preserves past knowledge, while larger  $\beta$  better maintains latest-task accuracy relative to the joint baseline.

Table 5: Mean accuracy and forgetting score across different  $(\alpha, \beta)$  settings and refinement epochs, evaluated with buffer sizes of 3K (tight budget) and 33K (standard in vision benchmarks).

Buffer	onoch	$ \begin{array}{c c} \textbf{Mean Accuracy} \\ \alpha = 1, \beta = 1 \ \alpha = 4, \beta = 1 \ \alpha = 1, \beta = 4 \end{array} $			Forgetting Score		
	epocii	$\alpha=1,\beta=1$	$\alpha \text{=} 4, \beta \text{=} 1$	$\alpha$ =1, $\beta$ =4	$\alpha=1,\beta=1$	$\alpha \text{=} 4, \beta \text{=} 1$	$\alpha \text{=} 1, \beta \text{=} 4$
3,000	8	46.67	47.77	45.39	49.05	47.06	52.87
	80	50.05	50.75	48.03	37.09	32.51	42.51
33,000	8	59.14	61	58.44	27.53	18.02	30.23
	80	60.26	61.35	59.22	21.13	16.72	25.4

Table 5 shows that mean accuracy changes little (about 3%) across  $(\alpha, \beta)$  settings under a 33,000 buffer and 8 epochs, but forgetting scores vary by up to 12% (highlighted in blue). For example, with  $\alpha=4,\beta=1$ , the forgetting scores for Tasks 2–5 are (11.74, 15.64, 18.60, 22.38), whereas with  $\alpha=1,\beta=4$  they rise to (19.11, 28.63, 32.98, 37.41). This indicates that larger  $\beta$  favors recent-task accuracy, while larger  $\alpha$  better retains earlier knowledge. Hence, tuning  $(\alpha,\beta)$  is essential not only for accuracy but also for managing the stability–plasticity trade-off in continual malware detection.

# 7 CONCLUSION

We presented TraMEL, a replay-based CL framework for malware traffic analysis. TraMEL mitigates catastrophic forgetting under class imbalance and temporal shifts, approaching joint-training performance while operating under strict memory budgets. Nonetheless, trade-offs remain: refinement may reduce current-task accuracy, and fixed buffers limit scalability. Future work should explore scalable backbones and exemplar-free methods to better handle highly imbalanced, dynamic environments.

#### REPRODUCIBILITY STATEMENT

We provide an anonymous GitHub repository containing the main source code used in our experiments. All datasets employed in this work (CICAndMal2017 and IoT-23) are publicly available, and we additionally release the preprocessing scripts to ensure consistent data preparation. The code includes a default seed, and we provide a hyperparameter as a default in the code to facilitate the reproduction of results with similar performance. While we do not provide strict hardware specifications, the implementation runs without specialized dependencies and has been tested under standard GPU environments. Overall, we aim to facilitate reproducibility by releasing both the code and preprocessing pipelines, enabling independent researchers to obtain comparable results.

#### ETHICS STATEMENT

This study aims to classify malware families in order to strengthen defenses against malicious attacks, with no intent of misuse. No human subjects are involved, and all datasets used (CICAnd-Mal2017, IoT-23) are publicly available.

## REFERENCES

- Suresh Kumar Amalapuram, Bheemarjuna Reddy Tamma, and Sumohana S Channappayya. SPI-DER: A semi-supervised continual learning-based network intrusion detection system. In <u>IEEE INFOCOM 2024-IEEE Conference on Computer Communications</u>, pp. 571–580. IEEE, 2024.
- Blake Anderson and David McGrew. Identifying encrypted malware traffic with contextual flow data. In <u>Proceedings of the 2016 ACM workshop on artificial intelligence and security</u>, pp. 35–46, 2016.
- Prashant Bhat, Bahram Zonooz, and Elahe Arani. Task-aware information routing from common representation space in lifelong learning. arXiv preprint arXiv:2302.11346, 2023a.
- Prashant Bhat, Bahram Zonooz, and Elahe Arani. Task-aware information routing from common representation space in lifelong learning. In <u>International Conference on Learning Representations (ICLR)</u>, 2023b.
- Sumohana Channappayya, Bheemarjuna Reddy Tamma, et al. Augmented memory replay-based continual learning approaches for network intrusion detection. <u>Advances in Neural Information</u> Processing Systems, 36:17156–17169, 2023.
- Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with A-GEM. In ICML, 2019.
- Yizheng Chen, Zhoujie Ding, and David Wagner. Continuous learning for android malware detection. In 32nd USENIX Security Symposium (USENIX Security 23), pp. 1127–1144, 2023.
- Xiyue Deng and Jelena Mirkovic. Polymorphic malware behavior through network trace analysis. In 2022 14th International Conference on COMmunication Systems & NETworkS (COMSNETS), pp. 138–146. IEEE, 2022.
- Sebastian Garcia, Agustin Parmisano, and Maria Jose Erquiaga. Iot-23: A labeled dataset with malicious and benign iot network traffic. (No Title), 2020.
- Xin Huang, Ashish Khetan, Milan Cvitkovic, and Zohar Karnin. Tabtransformer: Tabular data modeling using contextual embeddings. arXiv preprint arXiv:2012.06678, 2020.
- Alexander Küchler, Alessandro Mantovani, Yufei Han, Leyla Bilge, and Davide Balzarotti. Does every second count? time-based evolution of malware behavior in sandboxes. In Network and Distributed Systems Security Symposium. Internet Society, 2021.
- Arash Habibi Lashkari, Andi Fitriah A Kadir, Laya Taheri, and Ali A Ghorbani. Toward developing a systematic approach to generate benchmark android malware datasets and classification. In 2018 International Carnahan conference on security technology (ICCST), pp. 1–7. ieee, 2018.

- Zhizhong Li and Derek Hoiem. Learning without forgetting. <u>IEEE transactions on Pattern Analysis</u> and Machine Intelligence (TPAMI), 40(12):2935–2947, 2017.
- Yilin Lyu, Liyuan Wang, Xingxing Zhang, Zicheng Sun, Hang Su, Jun Zhu, and Liping Jing. et al. 2023). Advances in Neural Information Processing Systems, 36:25475–25494, 2023.
  - 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.
  - Yisroel Mirsky, Tomer Doitshman, Yuval Elovici, and Asaf Shabtai. Kitsune: An ensemble of autoencoders for online network intrusion detection. 2018.
  - Andreas Moser, Christopher Kruegel, and Engin Kirda. Limits of static analysis for malware detection. In Twenty-third annual computer security applications conference (ACSAC 2007), pp. 421–430. IEEE, 2007.
  - 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. <a href="arXiv"><u>arXiv</u></a> preprint arXiv:2501.01110, 2025.
  - Antonio Paya, Sergio Arroni, Vicente García-Díaz, and Alberto Gómez. Apollon: a robust defense system against adversarial machine learning attacks in intrusion detection systems. Computers & Security, 136:103546, 2024.
  - Mohammad Saidur Rahman, Scott E. Coull, and Matthew Wright. On the limitations of continual learning for malware classification. In <u>First Conference on Lifelong Learning Agents (CoLLAs)</u>, 2022.
  - Mohammad Saidur Rahman, Scott Coull, Qi Yu, and Matthew Wright. Madar: Efficient continual learning for malware analysis with distribution-aware replay. 2025. URL https://arxiv.org/abs/2502.05760.
  - Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. iCaRL: Incremental classifier and representation learning. In Conference on Computer Vision and Pattern Recognition (CVPR), 2017a.
  - Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In <u>Proceedings of the IEEE conference on Computer Vision and Pattern Recognition</u>, pp. 2001–2010, 2017b.
  - David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. Advances in neural information processing systems, 32, 2019.
  - Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. Advances in Neural Information Processing Systems (NeurIPS), 2017.
  - Robin Sommer and Vern Paxson. Outside the closed world: On using machine learning for network intrusion detection. In 2010 IEEE symposium on security and privacy, pp. 305–316. IEEE, 2010.
  - Tiezhu Sun, Nadia Daoudi, Weiguo Pian, Kisub Kim, Kevin Allix, Tegawende F Bissyande, and Jacques Klein. Temporal-incremental learning for android malware detection. <u>ACM Transactions on Software Engineering and Methodology</u>, 34(4):1–30, 2025.
  - Gido M van de Ven, Hava T Siegelmann, and Andreas S Tolias. Brain-inspired replay for continual learning with artificial neural networks. Nature Communications, 2020.
  - Xiaohu Xu, Xixi Zhang, Qianyun Zhang, Yu Wang, Bamidele Adebisi, Tomoaki Ohtsuki, Hikmet Sari, and Guan Gui. Advancing malware detection in network traffic with self-paced class incremental learning. <a href="IEEE Internet of Things Journal">IEEE Internet of Things Journal</a>, 11(12):21816–21826, 2024.
  - Xixi Zhang, Yu Wang, Tomoaki Ohtsuki, Guan Gui, Chau Yuen, Marco Di Renzo, and Hikmet Sari. Malware traffic classification via expandable class incremental learning with architecture search. IEEE Transactions on Information Forensics and Security, 2025.

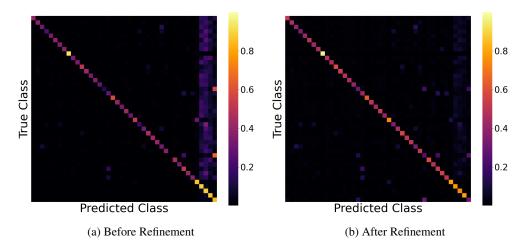


Figure 6: Class-level normalized confusion matrix of before and after refinement on CICAnd-Mal2017. Classes in recent task accuracy are over 0.8, while classes in earlier tasks have also been predicted as a last task.

Table 6: Task-wise accuracy and mean accuracy across tasks on the IoT23 dataset with different exemplar selection strategies.

Method	Task1	Task2	Task3	Task4	Task5	Mean
Random	98.38	90.54	79.94	69.91	73.39	82.43
C-Mean	98.38	88.46	79.74	65.19	72.18	80.79
$KM(N_k=5)$	98.38	83.00	72.60	66.60	71.33	78.382
$KM(N_k=100)$	98.38	90.01	80.06	66.95	72.04	81.49
$KM(N_k=300)$	98.38	88.98	79.71	70.20	73.03	82.06
$KM(N_k=600)$	98.38	90.78	80.40	69.28	75.08	82.78
$KM(N_k = 800)$	98.38	90.68	79.73	67.88	72.35	81.60
$KM(N_k=1,000)$	98.38	90.58	79.29	69.39	74.58	82.44

# A APPENDIX

# A.1 CONFUSION MATRIX

Using the CICAndMal2017 dataset, we trained over five tasks and analyzed class-wise normalized confusion matrices before and after refinement in the last task. As shown in Figure 6a, before refinement, many samples were misclassified into the latest task. After refinement 6b, the accuracy of all classes except the latest task improved, and overly predicting to the latest task is reduced.

# A.2 CLUSTER SIZE IN EXEMPLAR SELECTION

We evaluate exemplar selection strategies, random sampling (Random), class-mean selection as used in iCaRL (Rebuffi et al., 2017b) (C-Mean), and K-means clustering-based selection (KM). Experiments are conducted on the IoT23 dataset with a fixed buffer size of K=10,000, and the number of clusters  $N_k$  is varied from 5 to 1,000.

As shown in Table 6, TraMEL performance improves as  $N_k$  increases, peaking at around  $N_k = 600$ , and slightly degrades beyond that point. This trend suggests that increasing the number of clusters enhances class representation by promoting diversity in the selected exemplars. However, when  $N_k$  becomes too large, the number of samples per cluster becomes too small to capture intra-class variability, reducing the representativeness of the selected exemplars.

Overall, with  $K=10,000,\,K$ -means based selection yields more representative exemplars per class than random sampling or class-mean selection, leading to better overall performance.