

Shae McFadden^{1,2,3}, Myles Foley², Mario D'Onghia³, Chris Hicks²,
 Vasilios Mavroudis², Nicola Paoletti¹, Fabio Pierazzi³
¹King's College London, ²The Alan Turing Institute, ³University College London

Traditional Malware Detection

Thousands of new apps per day

Limited capacity for manual review

Concept Drift

ML Assumption: data is stationary

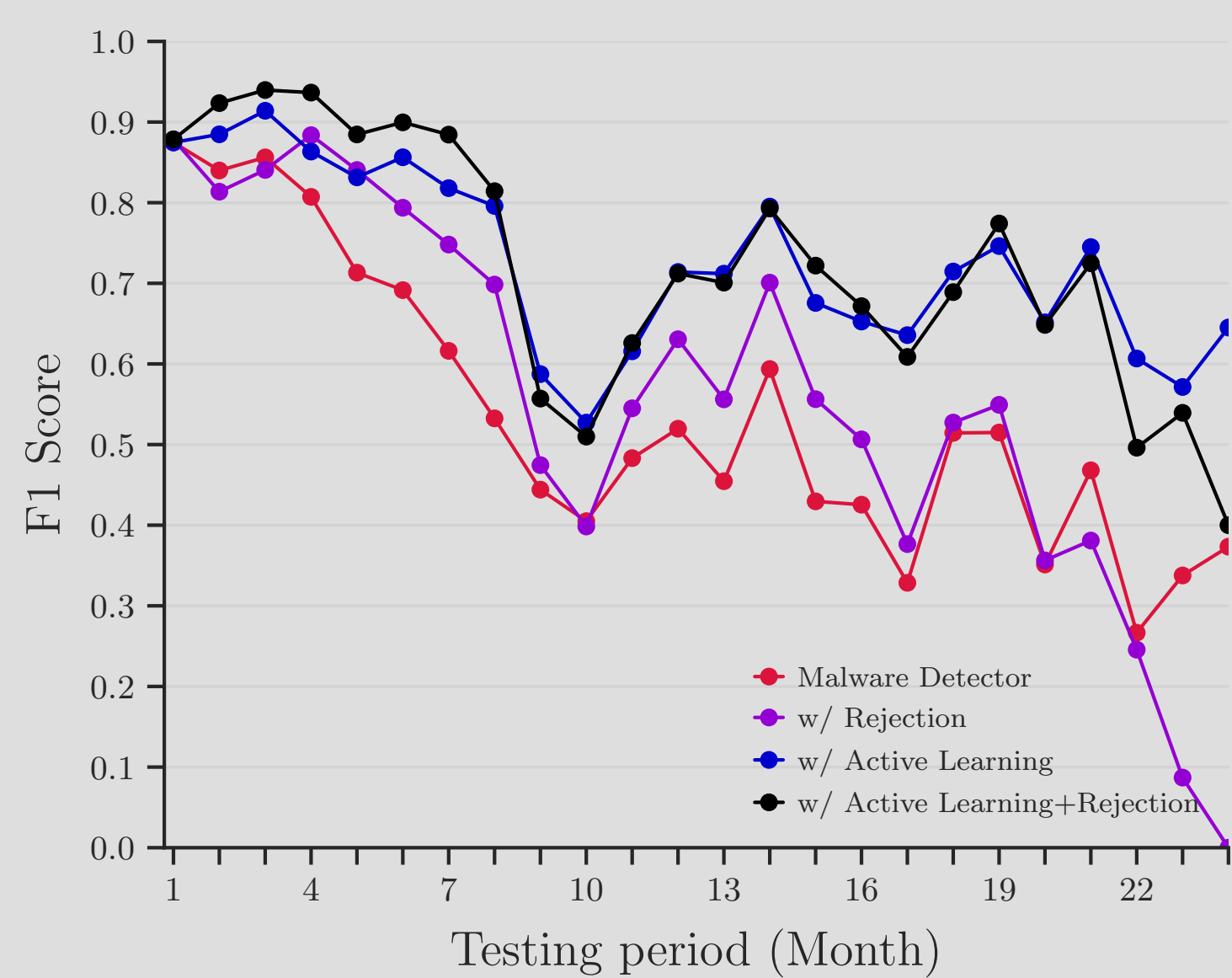
Reality: apps constantly evolve

Result: Performance degradation

Yesterday's training data becomes less relevant for today's threats.

Active Learning: Adapting the Detector to Drift

Selects an informative subset of new samples for retraining.



Feature Engineering

Malware/Goodware

Feature Extractor

Active Learning

Retraining

Manual Labelling

Selector

Classifier

Classified Samples

Rejection: Limiting the Impact of Drift

Selects samples at a high-risk of being misclassified to be quarantined.

Rejection

Rejector

Quarantined Samples

Key Observation

Existing approaches treat active learning, rejection, and detection independently

Intuition

Treat malware detection as a **unified** decision-making problem and use deep reinforcement learning

Rewards

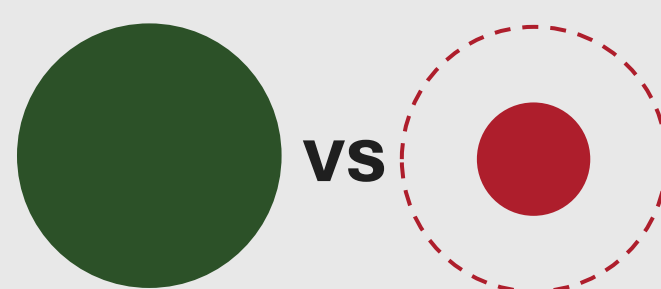
Accuracy

Provides the foundation
 +1 correct, -1 incorrect



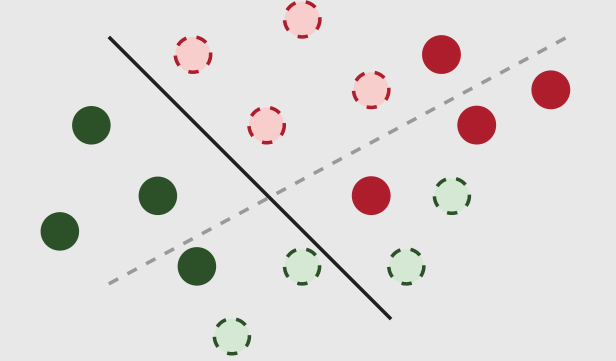
Class Imbalance

Upscales rewards for malware based on distribution (~10%)



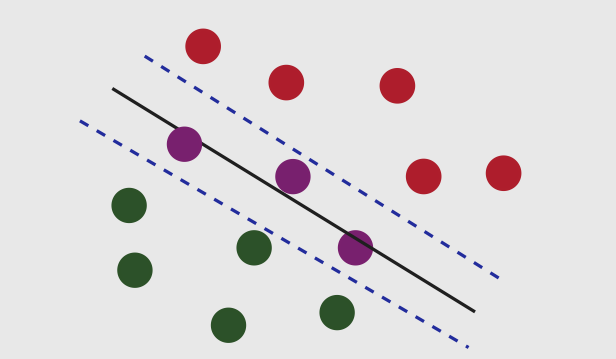
Temporal Robustness

Upscales rewards for samples based on temporal position



Rejection

Balances rewards for rejection relative to misclassification risk



Formulation (MD-MDP)

One-step MDP (Contextual Bandit)

Corrects spurious dependencies of prior work, ICMDP [Appl. Intell.'20]

Action Space

✓ Classify as Goodware
 ✗ Classify as Malware
 ? Reject → Active Learning

Experimental

Feature Spaces: Drebin (10,000D) and Ramda (379D)

Datasets: Hypercube (2021-2023) and Transcendent (2014-2018)

AMD Baselines: Drebin (SVM), DeepDrebin (MLP), and Ramda (MLP+VAE)

DRL Baselines: ICMDP and DCBs

MDP Comparison

Same CO policy architectures
 MD-MDP outperforms ICMDP

97% settings

45% significant

+1.94 ΔAUT

Classifier Comparison

Same AL and rejection budgets
 DRMD outperforms Baselines

90% settings

79% significant

+8.66 ΔAUT

Pipeline Comparison

Same AL and rejection budgets
 DRMD outperforms Baselines

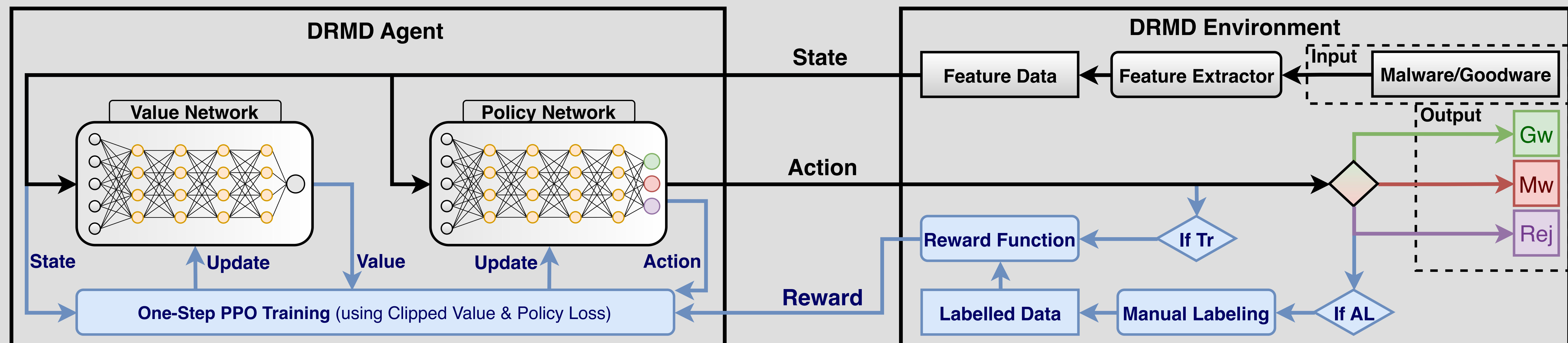
81% settings

68% significant

+10.90 ΔAUT

Takeaways

- 1) Adaptive decision-making, not just classification
- 2) One-step MDP formulation
- 3) Concept drift-aware DRL
- 4) Integration that matters
- 5) A starting point for future research



DRMD: Deep Reinforcement Learning for Malware Detection