

Malware Language Processing "MLP": Developing a new paradigm for malware analysis and classification using Machine Learning and Artificial Intelligence

Graduate Student : Solomon Sonya PI: Professor Dongyan Xu 2024 Department of Computer Science, Purdue University

1. Problem Statement

- Malware continues to increase in prevalence and sophistication.
- Traditional detection mechanisms including antivirus software fail to adequately detect new and varied malware.
- Artificial Intelligence and Machine Learning models (AI/ML) provide advanced capabilities that can be applied to the cybersecurity domain to enhance detection and classification of malware and malware families. Building a **robust and automated artificially intelligent malware analysis framework** and **producing new, standardized malware datasets for future classification in AI/ML** however, **are not trivial**.

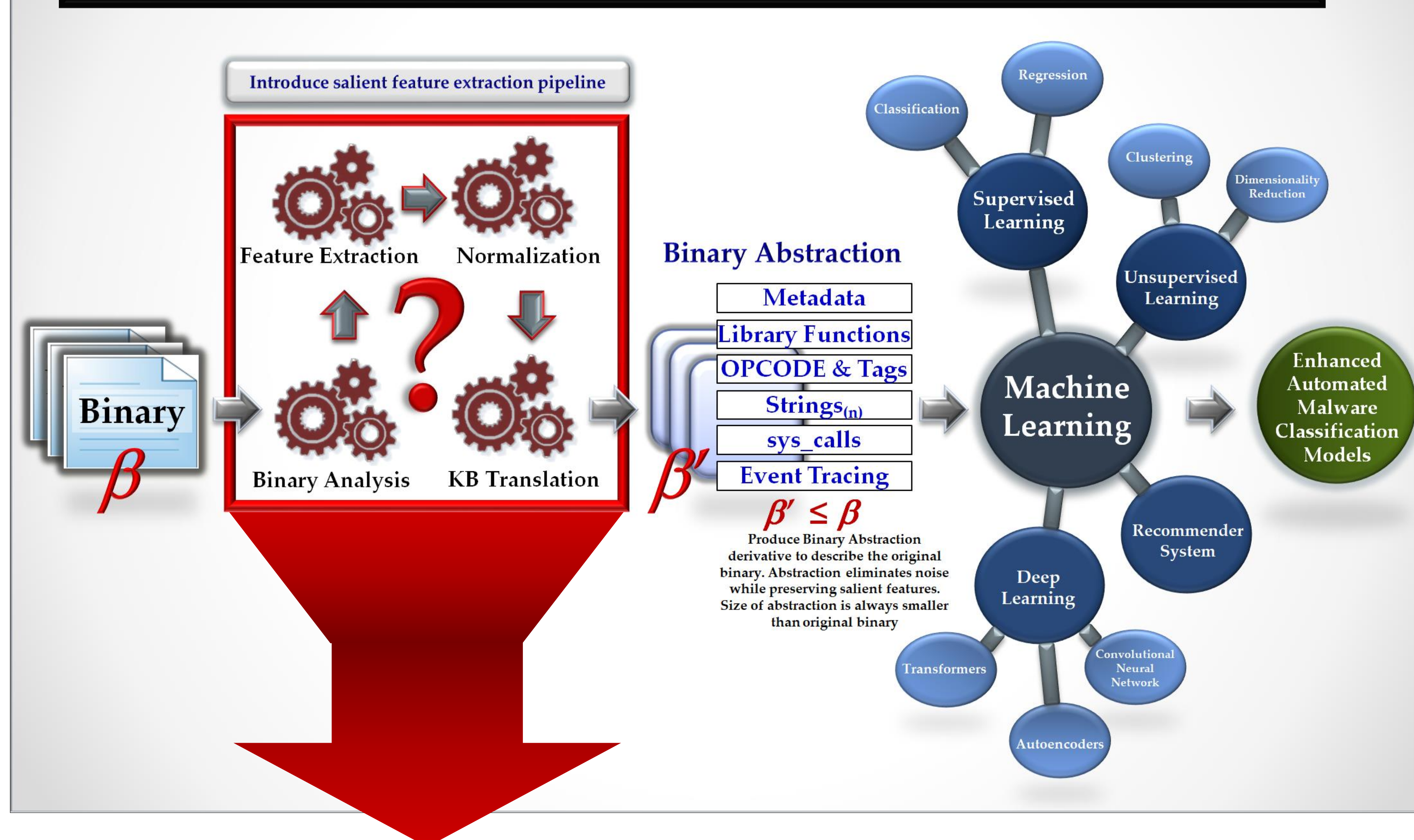
2. Research Objective

The principal objectives of this research are to deliver a new malware analysis framework, create a new ML heterogeneous classification model to analyze malware, produce uniform datasets for additional AI/ML analysis, and increase classification accuracy across complete and highly diverse malware corpora.

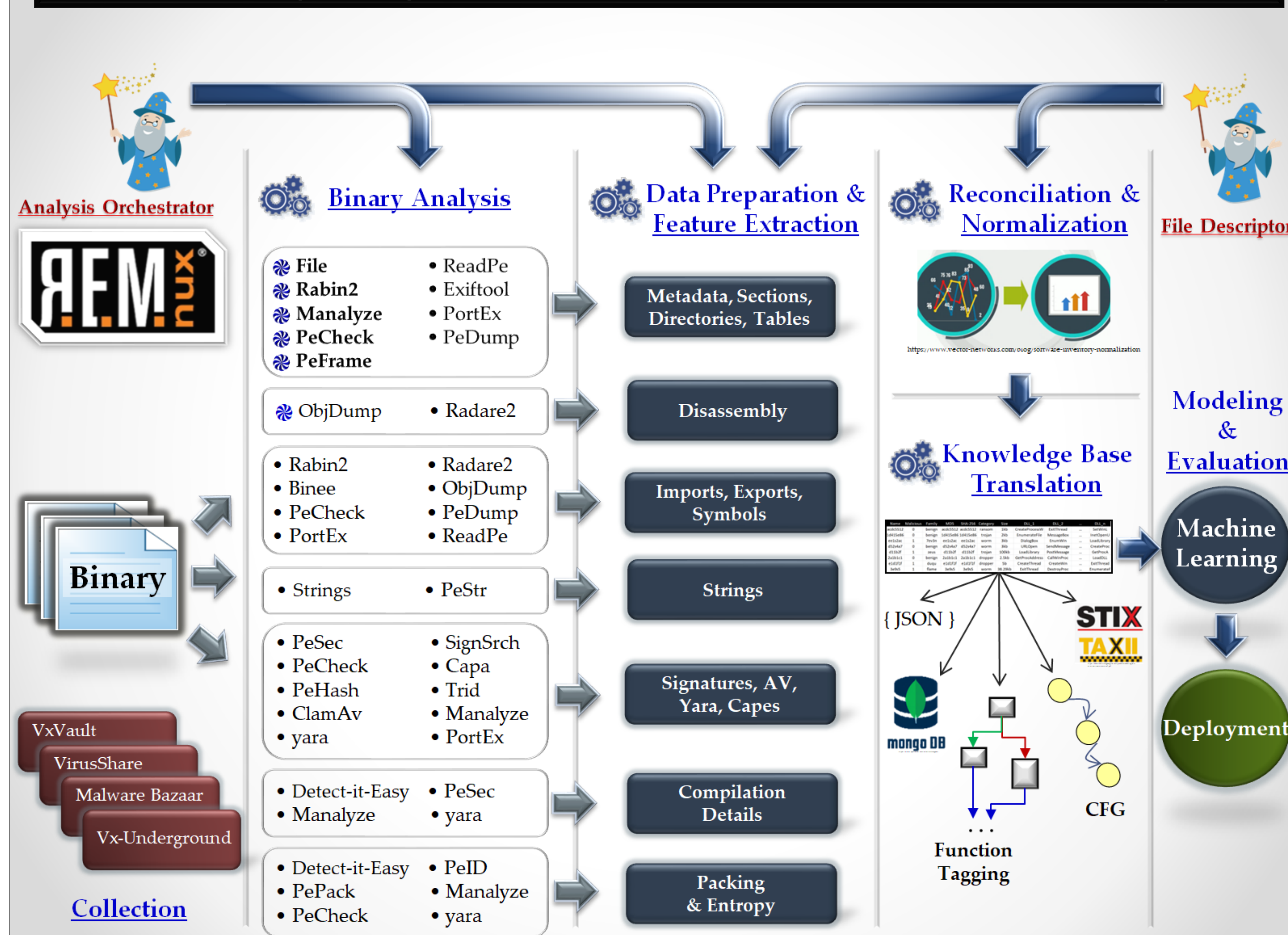
3. Research Contribution

- To our knowledge, this is the first research that:
- Releases a complete framework that automates malware analysis and produces new and standardized malware description datasets (ready for AI/ML analysis)
 - Develops a classification ensemble applied to malware analysis (named the Malware Ensemble Classification Facility) that improves traditional ML model selection.
 - **Releases trained ensemble classifier providing enhanced classification results of an entire 200+ gigabyte, malware family corpus consisting of 80K+ unique malware samples and 70+ malware families.**

4. Malware Analysis and Machine Learning Workflow

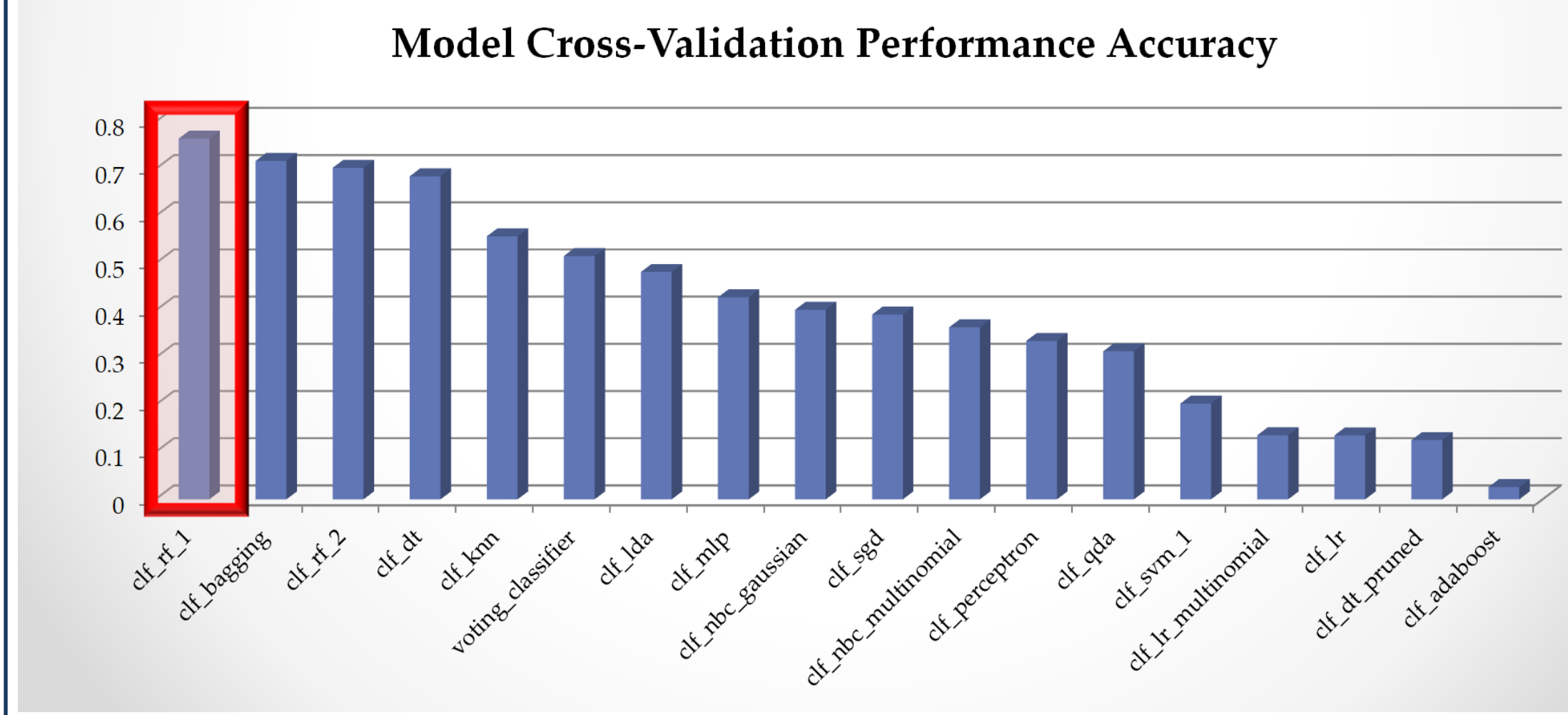


5. Static Binary Analysis Feature Extraction & Standardization Pipeline



6. Preliminary Analysis of Classification Models

18 ML Classifications Models Trained → Random Forest had best score



7. Classification Report of each Classifier Model

Ensemble classifier can improve results of overall classification performance

FYI: ClamAV detection rate: ~47%*

Model	7ev3n	9002Rat	Agent Testa	Aria Body	Bankshot	Bazarac Idoor	Bilster Loader	Cobalt Collector	Conti	CoreShell	Cryptolocker	Darkside	Dexia	Dyre
Best overall model	0.97	0.97	0.97	0.97	0.97	0.97	0.97	0.97	0.97	0.97	0.97	0.97	0.97	0.97
Superior Model for Specific Class	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99

8. Conclusions & Future Work

A new framework has been developed that currently automates static malware analysis, performs feature extraction, and standardizes malware dataset production for AI/ML analysis. 18 ML classification models were trained on a sampling (3100 binaries representing 111 unique malware families) of our entire malware corpus of 80K+ binaries. Random Forests produced the best performance, however, we identify models with superior classification performance across discrete malware families. Future work includes enriching dataset production and AI/ML analysis and completing development of our new ensemble classifier named the **Malware Ensemble Classification Facility**.