# Automatic Program Protection Using Fuzzing Driven Classifiers
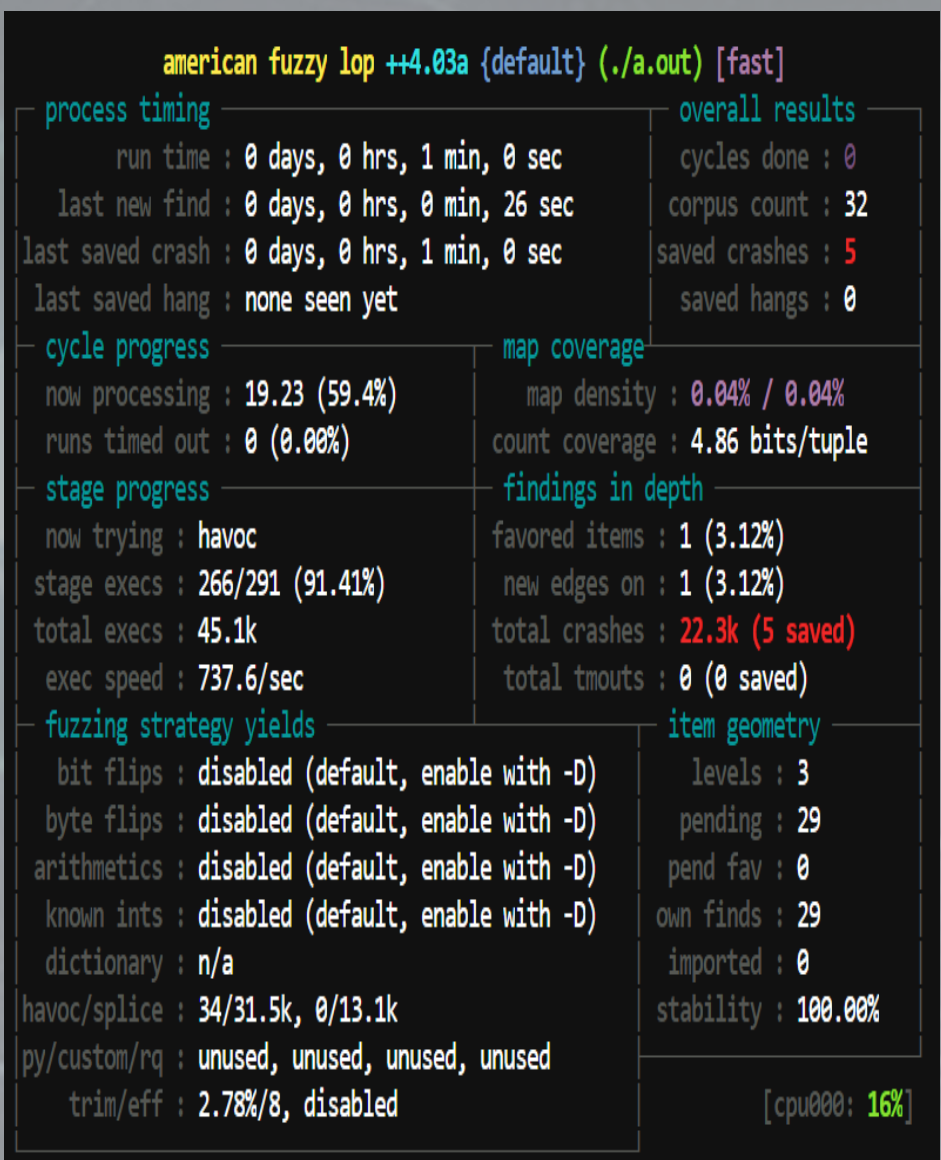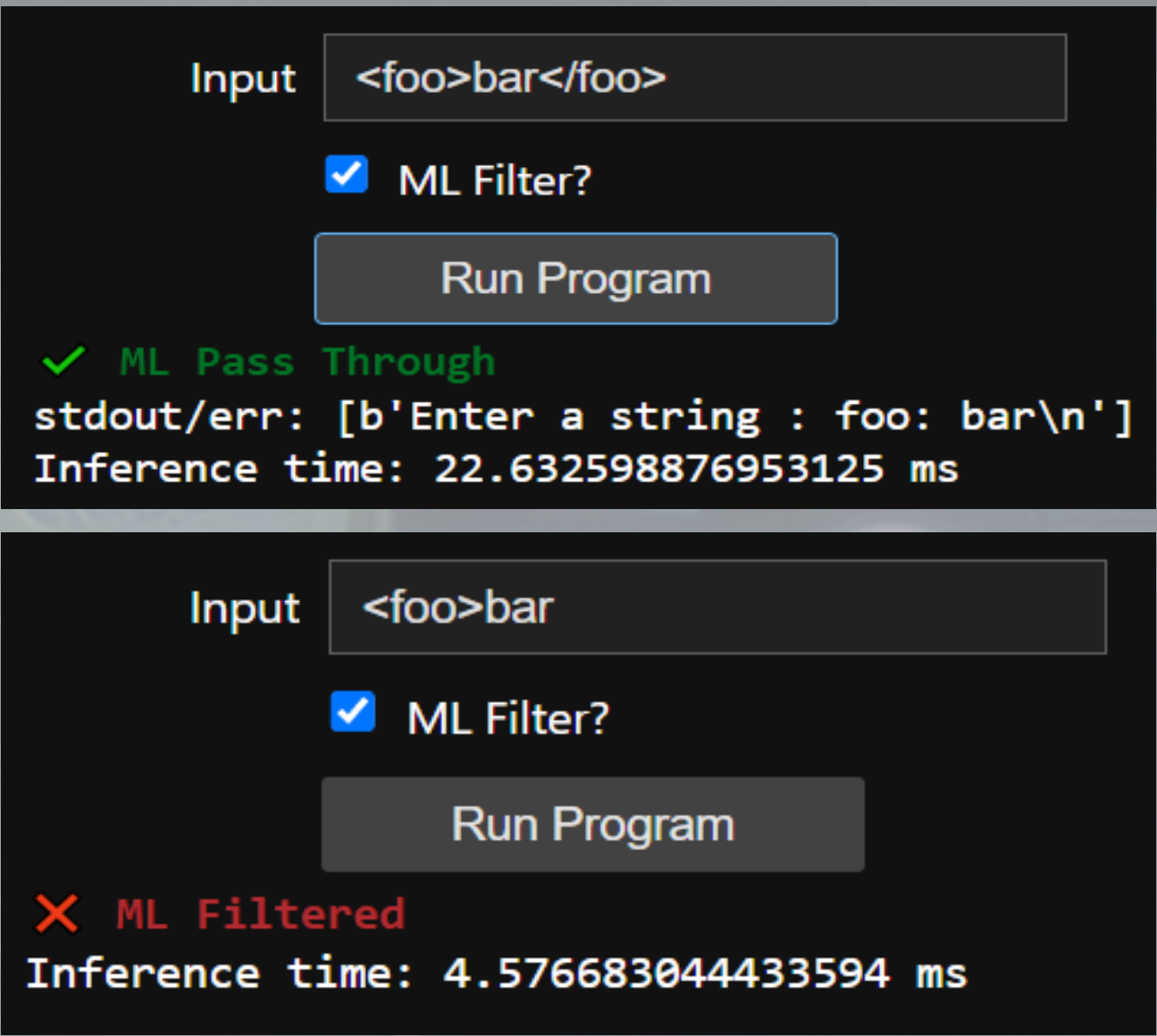
## Introduction

### Goal
- Classify program inputs to avoid crashing states that may lead to vulnerabilities
- Using fuzzer driven ML models to create automatic program protection filters
- Leveraging novel techniques in generating high-quality data for attention-based architectures

### Address Limitations in Past Work
- Generates large (500k+) training datasets
- Usage for protection and not fuzzing guidance
- Models capable of extracting syntactic features
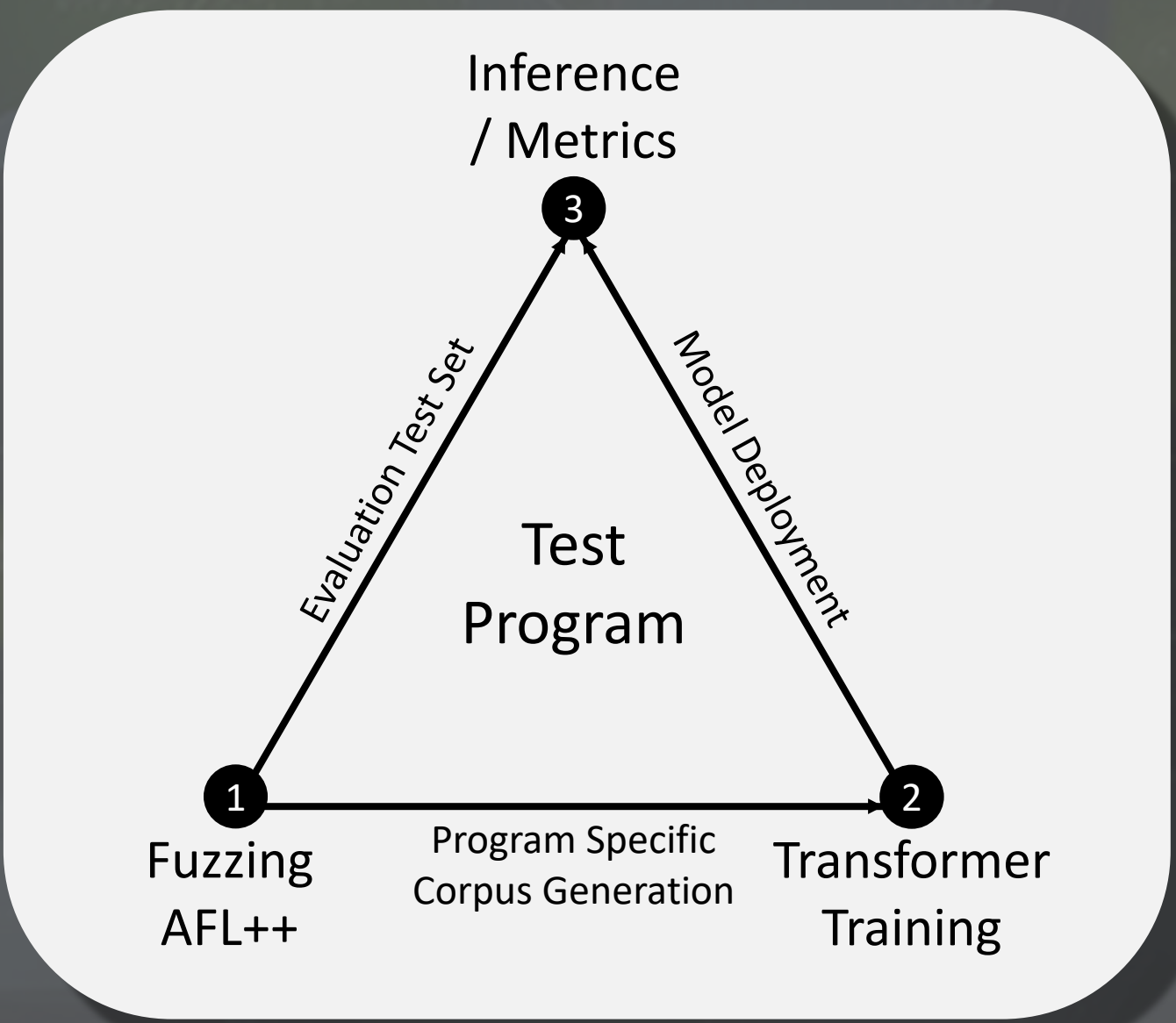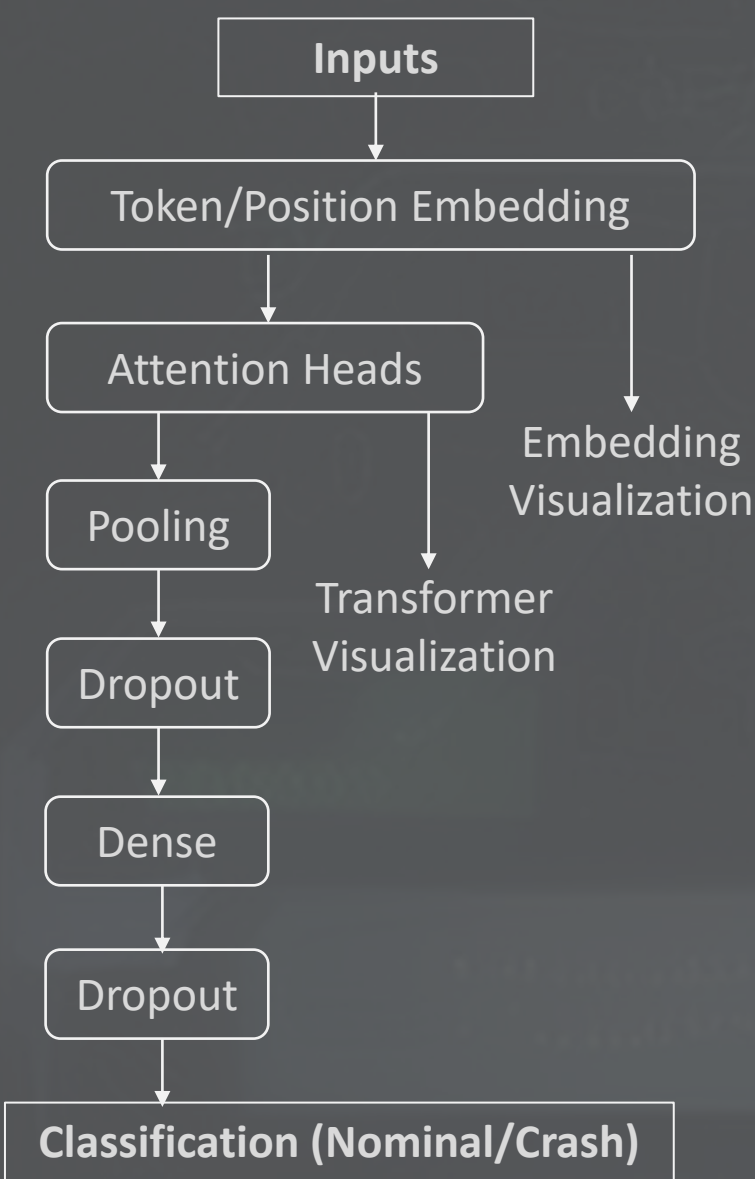- Analysis of corpus coverage and layer contributions



**Automatic Supervised Learning Corpus**



**Provide Tailored Program Protection via Classifier**

## Transformer Architecture



## Experiment

### Approach
1. Modified fuzzer to generate training data for supervised learning
2. Binary classifier model trained with attention heads to learn nominal vs. crashing
3. Inference for collecting metrics and providing tailored application protection

### Metrics
- Ablation study: Statistics collected for reduced models: dense only, positional, transformer
- F1/AUC Scores: Measure performance of classifier with false positives/negatives
- PaCMAP: Qualitative analysis of manifold coverage for both datasets and activations

## Findings

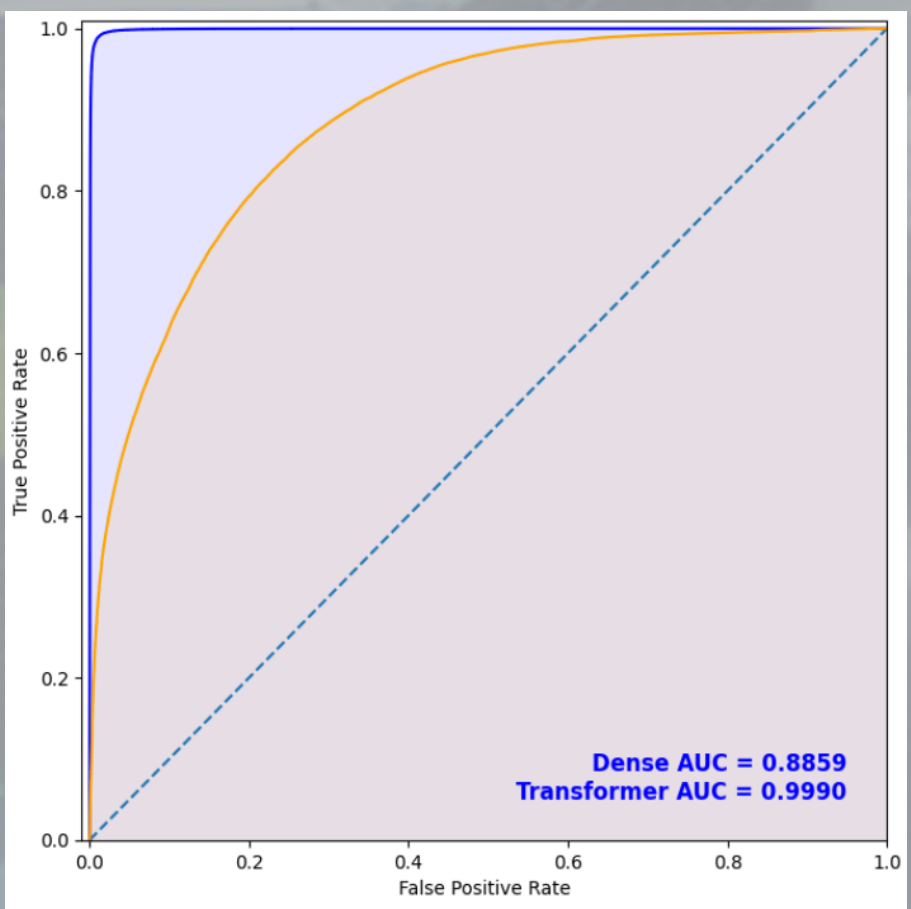| | Dense | | Position | | Attention | |
|---|---|---|---|---|---|---|
| | F1 | AUC | F1 | AUC | F1 | AUC |
| Test | 0.8044 | 0.8666 | 0.9695 | 0.9880 | 0.9911 | 0.9993 |
| Fuzzgoat | 0.9527 | 0.9851 | 0.9355 | 0.9746 | 0.9906 | 0.9994 |
| XML CVE | 0.9744 | 0.9747 | 0.9860 | 0.9964 | 0.9992 | 1.0000 |

### Quantitative
- Transformer model more performant in all cases
- Sequence information needed for inputs with heavy overlap, as seen in Test Program
- Able to correctly identify CVE inducing inputs
- Dense model achieves similar performance when special characters are present
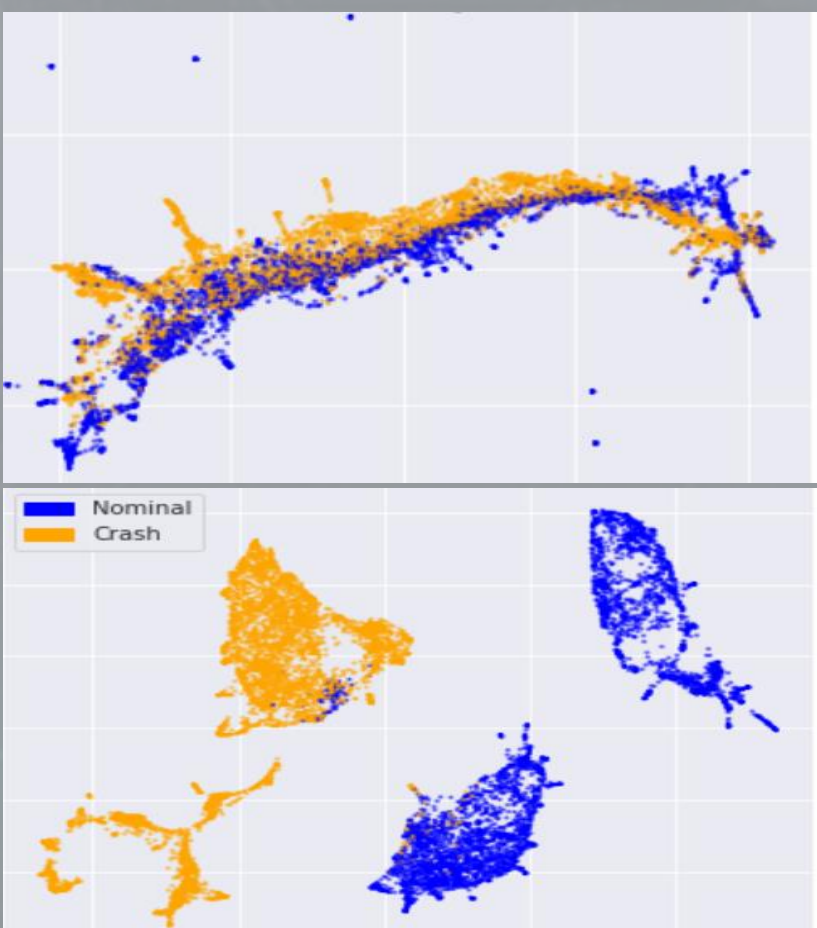
### Qualitative
- Good solution-space coverage from fuzzer, acts as a form of dataset augmentation
- Extra structure added by transformer to create distinct groupings on learned manifold



**Test Classifier Scores**



**Learned Dataset Structure**